# Robust Distillation for Worst-class Performance:
# On the Interplay Between Teacher and Student Objectives
# (Supplementary Material)

**Serena Wang**[1,2]     **Harikrishna Narasimhan**[1]     **Yichen Zhou**[1]     **Sara Hooker**[3]     **Michal Lukasik**[1]

**Aditya Krishna Menon**[1]

[1]Google Research, Mountain View, California and New York, New York, USA,
[2]University of California, Berkeley, Berkeley, California, USA,
[3]Cohere For AI, Palo Alto, California, USA

## A   PROOFS

### A.1   PROOF OF THEOREM 1

(i) The first result follows from the fact that the cross-entropy loss is a proper composite loss [Williamson et al., 2016] with the softmax function as the associated (inverse) link function.

(ii) For a proof of the second result, please see Menon et al. [2021b].

(iii) Below, we provide a proof for the third result.

The minimization of the robust objective in (3) over $f$ can be re-written as a min-max optimization problem:

$$\min_{f:\mathcal{X}\to\mathbb{R}^m} L^{\text{rob}}(f) = \min_{f:\mathcal{X}\to\mathbb{R}^m} \max_{\lambda\in\Delta_m} \underbrace{\sum_{y=1}^{m} \frac{\lambda_y}{\pi_y} \mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right]}_{\omega(\lambda,f)}. \tag{1}$$

The min-max objective $\omega(\lambda, f)$ is clearly linear in $\lambda$ (for fixed $f$) and with $\ell$ chosen to be the cross-entropy loss, is convex in $f$ (for fixed $\lambda$), i.e., $\omega(\lambda, \kappa f_1 + (1-\kappa)f_2) \leq \kappa\omega(\lambda, f_1) + (1-\kappa)\omega(\lambda, f_2)$, $\forall f_1, f_2 : \mathcal{X}\to\mathbb{R}^m, \kappa\in[0,1]$. Furthermore, $\Delta_m$ is a convex compact set, while the domain of $f$ is convex. It follows from Sion's minimax theorem [Sion, 1958] that:

$$\min_{f:\mathcal{X}\to\mathbb{R}^m} \max_{\lambda\in\Delta_m} \omega(\lambda, f) = \max_{\lambda\in\Delta_m} \min_{f:\mathcal{X}\to\mathbb{R}^m} \omega(\lambda, f). \tag{2}$$

Let $(\lambda^*, f^*)$ be such that:

$$\lambda^* \in \operatorname*{argmax}_{\lambda\in\Delta_m} \min_{f:\mathcal{X}\to\mathbb{R}^m} \omega(\lambda, f); \quad f^* \in \operatorname*{argmin}_{f:\mathcal{X}\to\mathbb{R}^m} \max_{\lambda\in\Delta_m} \omega(\lambda, f),$$

Such a $\lambda^*$ exists for the following reason: for any fixed $\lambda\in\Delta_m$, owing to the use of the cross-entropy loss, a minimizer over always exists for $\omega(\lambda, f)$, and is given by $f_y(x) = \log\left(\frac{\lambda_y}{\pi_y}\eta_y(x)\right) + C$, for some $C\in\mathbb{R}$; therefore $\min_{f:\mathcal{X}\to\mathbb{R}^m}\omega(\lambda, f)$ is bounded above for any $\lambda$, and $\Delta_m$ being compact set gives us there exits a maximizer $\lambda^*$ over this set. Similarly, such an $f^*$ exists for the following reason: the objective $\max_{\lambda\in\Delta_m}\omega(\lambda, f)$ takes a bounded value when $f = \eta$, and any minimizer of $\max_{\lambda\in\Delta_m}\omega(\lambda, f)$ yields a value below that; because $\omega(\lambda, f) \geq 0$ and is convex in $f$, the minimizer $f^*$ exits.

We then have from (2):

$$\begin{aligned}
\omega(\lambda^*, f^*) &\leq \max_{\lambda\in\Delta_m} \omega(\lambda, f^*) \\
&= \min_{f:\mathcal{X}\to\mathbb{R}^m} \max_{\lambda\in\Delta_m} \omega(\lambda, f) = \max_{\lambda\in\Delta_m} \min_{f:\mathcal{X}\to\mathbb{R}^m} \omega(\lambda, f) \\
&= \min_{f:\mathcal{X}\to\mathbb{R}^m} \omega(\lambda^*, f) \leq \omega(\lambda^*, f^*),
\end{aligned}$$

which tells us that there exists $(\lambda^*, f^*)$ is a saddle-point for (1), i.e.,

$$\omega(\lambda^*, f^*) = \max_{\lambda \in \Delta_m} \omega(\lambda, f^*) = \min_{f:\mathcal{X} \to \mathbb{R}^m} \omega(\lambda^*, f).$$

Consequently, we have:

$$L^{\text{rob}}(f^*) = \max_{\lambda \in \Delta_m} \omega(\lambda, f^*) = \min_{f:\mathcal{X} \to \mathbb{R}^m} \max_{\lambda \in \Delta_m} \omega(\lambda, f) = \min_{f:\mathcal{X} \to \mathbb{R}^m} L^{\text{rob}}(f).$$

We thus have that $f^*$ is a minimizer of $L^{\text{rob}}(f)$. Furthermore, because $f^*$ is also a minimizer of $\omega(\lambda^*, f)$ over $f$, i.e.,

$$f^* \in \operatorname*{argmin}_{f:\mathcal{X} \to \mathbb{R}^m} \sum_{y=1}^{m} \frac{\lambda_y^*}{\pi_y} \mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right],$$

it follows that:

$$\operatorname{softmax}_y(f^*(x)) \propto \frac{\lambda_y^*}{\pi_y} \eta_y(x).$$

(iv) For the fourth result, we expand the traded-off objective, and re-write it as:

$$
\begin{aligned}
L^{\text{tdf}}(f) &= (1 - \alpha)L^{\text{bal}}(f) + \alpha L^{\text{rob}}(f) \\
&= (1 - \alpha)\frac{1}{m}\sum_{y=1}^{m}\frac{1}{\pi_y}\mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right] + \alpha \max_{\lambda \in \Delta_m}\sum_{y=1}^{m}\frac{\lambda_y}{\pi_y}\mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right] \\
&= \max_{\lambda \in \Delta_m}\underbrace{\sum_{y=1}^{m}\left((1-\alpha)\frac{1}{m} + \alpha\lambda_y\right)\frac{1}{\pi_y}\mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right]}_{\omega(\lambda, f)}.
\end{aligned}
$$

For a fixed $\lambda$, $\omega(\lambda, f)$ is convex in $f$ (as the loss $\ell$ is the cross-entropy loss), and for a fixed $f$, $\omega(\lambda, f)$ is linear in $\lambda$. Following the same steps as the proof of (iii), we have that there exists $(\lambda^*, f^*)$ such that

$$L^{\text{tdf}}(f^*) = \max_{\lambda \in \Delta_m} \omega(\lambda, f^*) = \min_{f:\mathcal{X} \to \mathbb{R}^m} L^{\text{tdf}}(f),$$

and

$$f^* \in \operatorname*{argmin}_{f:\mathcal{X} \to \mathbb{R}^m} \sum_{y=1}^{m}\left((1-\alpha)\frac{1}{m} + \alpha\lambda_y^*\right)\frac{1}{\pi_y}\mathbb{E}\left[\eta_y(X)\,\ell(y, f(X))\right],$$

which, owing to the properties of the cross-entropy loss, then gives us the desired form for $f^*$.

### A.2   PROOF OF THEOREM 2

*Proof.* Expanding the left-hand side, we have:

$$
\begin{aligned}
|\hat{L}^{\text{rob-d}}(f) - L^{\text{rob}}(f)| &\leq |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f) + L^{\text{rob-d}}(f) - L^{\text{rob}}(f)| \\
&\leq |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f)| + |L^{\text{rob-d}}(f) - L^{\text{rob}}(f)| \\
&= |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f)| + \left|\max_{y \in [m]}\frac{\mathbb{E}_x\left[p_y^t(x)\,\ell(y, f(x))\right]}{\mathbb{E}_x\left[p_y^t(x)\right]} - \max_{y \in [m]}\frac{\mathbb{E}_x\left[\eta_y(x)\,\ell(y, f(x))\right]}{\pi_y}\right| \\
&\leq |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f)| + \max_{y \in [m]}\left|\frac{\mathbb{E}_x\left[p_y^t(x)\,\ell(y, f(x))\right]}{\mathbb{E}_x\left[p_y^t(x)\right]} - \frac{\mathbb{E}_x\left[\eta_y(x)\,\ell(y, f(x))\right]}{\pi_y}\right| \\
&\leq |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f)| + B\max_{y \in [m]}\mathbb{E}_x\left[\left|\frac{p_y^t(x)}{\mathbb{E}_x\left[p_y^t(x)\right]} - \frac{\eta_y(x)}{\pi_y}\right|\ell(y, f(x))\right] \\
&\leq |\hat{L}^{\text{rob-d}}(f) - L^{\text{rob-d}}(f)| + B\max_{y \in [m]}\mathbb{E}_x\left[\left|\frac{p_y^t(x)}{\mathbb{E}_x\left[p_y^t(x)\right]} - \frac{\eta_y(x)}{\pi_y}\right|\right],
\end{aligned}
$$

where the second-last step uses Jensen's inequality and the fact that $\ell(y, f(x)) \geq 0$, and the last step uses the fact that $\ell(y, f(x)) \leq B$.

Further expanding the first term,

$$\left| \hat{L}^{\text{rob-d}}(f) - L^{\text{rob}}(f) \right| \leq \left| \max_{y \in [m]} \phi_y(f) - \max_{y \in [m]} \hat{\phi}_y(f) \right| + B \max_{y \in [m]} \mathbb{E}_x \left[ \left| \frac{p_y^t(x)}{\mathbb{E}_x \left[ p_y^t(x) \right]} - \frac{\eta_y(x)}{\pi_y} \right| \right]$$

$$\leq \max_{y \in [m]} \left| \phi_y(f) - \hat{\phi}_y(f) \right| + B \max_{y \in [m]} \mathbb{E}_x \left[ \left| \frac{p_y^t(x)}{\mathbb{E}_x \left[ p_y^t(x) \right]} - \frac{\eta_y(x)}{\pi_y} \right| \right],$$

as desired. $\qquad \square$

## A.3 CALIBRATION OF MARGIN-BASED LOSS

To show that minimizer of the margin-based objective in (9) also minimizes the balanced objective in (6), we state the following general result:

**Lemma 1.** *Suppose $p^t \in \mathcal{F}$ and $\mathcal{F}$ is closed under linear transformations. Let*

$$\hat{f} \in \operatorname*{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}^{\text{mar}} \left( p^t(x_i), f(x_i); \mathbf{c} \right) \tag{3}$$

*for some cost vector $\mathbf{c} \in \mathbb{R}_+^m$. Then:*

$$\hat{f}_y(x_i) = \log \left( c_y p_y^t(x_i) \right) + C_i, \quad \forall i \in [n],$$

*for some example-specific constant constants $C_i \in \mathbb{R}, \forall i \in [n]$. Furthermore, for any assignment of example weights of $w \in \mathbb{R}_+^n$, $\hat{f}$ is also the minimizer of the weighted objective:*

$$\hat{f} \in \operatorname*{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} w_i \sum_{y=1}^{m} c_y \, p_y^t(x_i) \, \ell \left( y, f(x_i) \right). \tag{4}$$

*Proof.* Following Menon et al. [2021b] (e.g. proof of Theorem 1), we have that for class probabilities $\mathbf{p} \in \Delta_m$ and costs $\mathbf{c} \in \mathbb{R}_+^m$, the margin-based loss in (9)

$$\mathcal{L}^{\text{mar}} \left( \mathbf{p}, \mathbf{f}; \mathbf{c} \right) = \frac{1}{m} \sum_{y \in [m]} p_y \log \left( 1 + \sum_{j \neq y} \exp \left( \log(c_y/c_j) - (f_y - f_j) \right) \right).$$

is minimized by:

$$f_y^* = \log \left( c_y p_y \right) + C,$$

for any $C > 0$. To see why this is true, note that the above loss can be equivalently written as:

$$\mathcal{L}^{\text{mar}} \left( \mathbf{p}, \mathbf{f}; \mathbf{c} \right) = -\frac{1}{m} \sum_{y \in [m]} p_y \log \left( \frac{\exp \left( f_y - \log(c_y) \right)}{\sum_{j=1}^{m} \exp \left( f_j - \log(c_j) \right)} \right).$$

This the same as the softmax cross-entropy loss with adjustments made to the logits, the minimizer for which is of the form:

$$f_y^* - \log(c_y) = \log \left( p_y \right) + C \quad \text{or} \quad f_y^* = \log \left( c_y p_y \right) + C.$$

It follows that any minimizer $\hat{f}$ of the average margin-based loss in (3) over sample $S$, would do so point-wise, and therefore

$$\hat{f}_y(x_i) = \log \left( c_y p_y^t(x_i) \right) + C_i, \quad \forall i \in [n],$$

for some example-specific constant constants $C_i \in \mathbb{R}, \forall i \in [n]$.

To prove the second part, we note that for the minimizer $\hat{f}$ to also minimize the weighted objective:

$$\frac{1}{n}\sum_{i=1}^{n} w_i \sum_{y=1}^{m} c_y \, p_y^t(x_i) \, \ell\left(y, f(x_i)\right),$$

it would also have to do so point-wise for each $i \in [m]$, and so as long the weights $w_i$ are non-negative, it suffices that

$$\hat{f}(x_i) \in \underset{\mathbf{f} \in \mathbb{R}^m}{\operatorname{argmin}} \sum_{y=1}^{m} c_y \, p_y^t(x_i) \, \ell\left(y, f(x_i)\right).$$

This is indeed the case when $\ell$ is the softmax cross-entropy loss, where the point-wise minimizer for each $i \in [m]$ would be of the form $\operatorname{softmax}_y(f(x)) = c_y p_y^t(x)$, which is satisfied by $\hat{f}$. $\qquad \square$

A similar result also holds in the population limit, when (3) and (4) are computed in expectation, and the per-example weighting in (4) is replaced by an arbitrary weighting function $w(x) \in \mathbb{R}_+$. Any scorer of the following form would then minimize both objectives:

$$\hat{f}_y(x) = \log\left(c_y p_y^t(x)\right) + C(x), \quad \forall x \in \mathcal{X},$$

where $C(x)$ is some example-specific constant.

## A.4 PROOF OF PROPOSITION 3

**Proposition** (Restated). *Suppose $p^t \in \mathcal{F}$ and $\mathcal{F}$ is closed under linear transformations. Then the final scoring function $\bar{f}^s(x) = \frac{1}{K}\sum_{k=1}^{K} f^k(x)$ output by Algorithm 1 is of the form:*

$$\operatorname{softmax}_j(\bar{f}^s(x)) \propto \bar{\lambda}_j p_j^t(x), \quad \forall j \in [m], \ \forall (x, y) \in S,$$

*where $\bar{\lambda}_y = \left(\prod_{k=1}^{K} \lambda_y^k / \pi_y^t\right)^{1/K}$.*

*Proof.* The proof follows from Lemma 1 with the costs $\mathbf{c}$ set to $\lambda^k / \pi^t$ for each iteration $k$. The lemma tells us that each $f^k$ is of the form:

$$f^k(x') = \log\left(\frac{\lambda_y^k}{\pi_y^t} p_y^t(x')\right) + C(x'), \quad \forall (x', y') \in S,$$

for some example-specific constant $C(x') \in \mathbb{R}$. Consequently, we have that:

$$\bar{f}_y^s(x') = \log(\bar{\lambda}_y p_y^t(x')) + \bar{C}(x'), \quad \forall (x', y') \in S,$$

where $\bar{\lambda}_y = \left(\prod_{k=1}^{K} \lambda_y^k / \pi_y^t\right)^{1/K}$ and $\bar{C}(x') \in \mathbb{R}$. Applying a softmax to $\bar{f}^s$ results in the desired form. $\qquad \square$

## A.5 PROOF OF THEOREM 4

**Theorem** (Restated). *Suppose $p^t \in \mathcal{F}$ and $\mathcal{F}$ is closed under linear transformations. Suppose $\ell$ is the softmax cross-entropy loss $\ell^{\mathrm{xent}}$, $\ell(y, z) \leq B$ and $\max_{y \in [m]} \frac{1}{\pi_y^t} \leq Z$, for some $B, Z > 0$. Furthermore, suppose for any $\delta \in (0, 1)$, the following bound holds on the estimation error in Theorem 2: with probability at least $1 - \delta$ (over draw of $S \sim D^n$), for all $f \in \mathcal{F}$,*

$$\max_{y \in [m]} \left|\phi_y(f) - \hat{\phi}_y(f)\right| \leq \Delta(n, \delta),$$

for some $\Delta(n,\delta) \in \mathbb{R}_+$ *that is increasing in* $1/\delta$, *and goes to 0 as* $n \to \infty$. *Fix* $\delta \in (0,1)$. *Then when the step size* $\gamma = \frac{1}{2BZ}\sqrt{\frac{\log(m)}{K}}$ *and* $n^{\text{val}} \geq 8Z\log(2m/\delta)$, *with probability at least* $1 - \delta$ *(over draw of* $S \sim D^n$ *and* $S^{\text{val}} \sim D^{n^{\text{val}}}$)

$$L^{\text{rob}}(\bar{f}^s) \leq \min_{f \in \mathcal{F}} L^{\text{rob}}(f) + \underbrace{2B \max_{y \in [m]} \mathbb{E}_x \left[\left|\frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right]}_{\text{Approximation error}}$$

$$+ \underbrace{2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)}_{\text{Estimation error}} + \underbrace{4BZ\sqrt{\frac{\log(m)}{K}}}_{\text{EG convergence}}.$$

Before proceeding to the proof, we will find it useful to define:

$$\hat{\phi}_y^{\text{val}}(f^s) = \frac{1}{\hat{\pi}_y^{t,\text{val}}}\frac{1}{n^{\text{val}}} \sum_{(x',y') \in S^{\text{val}}} p_y^t(x')\, \ell(y, f^s(x')).$$

We then state a useful lemma.

**Lemma 2.** *Suppose the conditions in Theorem 4 hold. Then with probability* $\leq 1 - \delta$ *(over draw of* $S \sim D^n$ *and* $S^{\text{val}} \sim D^{n^{\text{val}}}$), *at each iteration* $k$,

$$\sum_{y=1}^m \lambda_y^{k+1} \phi_y(f^{k+1}) - \min_{f \in \mathcal{F}} \sum_{y=1}^m \lambda_y^{k+1}\phi_y(f) \leq 2\Delta(n,\delta);$$

*and for any* $\lambda \in \Delta_m$:

$$\left|\sum_{y=1}^m \lambda_y \hat{\phi}_y^{\text{val}}(f^{k+1}) - \sum_{y=1}^m \lambda_y \phi_y(f^{k+1})\right| \leq \Delta(n^{\text{val}}, \delta).$$

*Proof.* We first note that by applying Lemma 1 with $w_i = 1, \forall i$, we have that $f^{k+1}$ is the minimizer of $\sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f)$ over all $f \in \mathcal{F}$, and therefore:

$$\sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f^{k+1}) \leq \sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f), \ \forall f \in \mathcal{F}. \tag{5}$$

Further, for a fixed iteration $k$, let us denote $\tilde{f} \in \operatorname*{argmin}_{f \in \mathcal{F}} \sum_{y=1}^m \lambda_y^{k+1}\phi_y(f)$. Then for the first part, we have:

$$\sum_{y=1}^m \lambda_y^{k+1}\phi_y(f^{k+1}) - \sum_{y=1}^m \lambda_y^{k+1}\phi_y(\tilde{f})$$

$$\leq \sum_{y=1}^m \lambda_y^{k+1}\phi_y(f^{k+1}) - \sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f^{k+1}) + \sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f^{k+1}) - \sum_{y=1}^m \lambda_y^{k+1}\phi_y(\tilde{f})$$

$$\leq \sum_{y=1}^m \lambda_y^{k+1}\phi_y(f^{k+1}) - \sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f^{k+1}) + \sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(\tilde{f}) - \sum_{y=1}^m \lambda_y^{k+1}\phi_y(\tilde{f})$$

$$\leq 2 \sup_{f \in \mathcal{F}} \left|\sum_{y=1}^m \lambda_y^{k+1}\hat{\phi}_y(f) - \sum_{y=1}^m \lambda_y^{k+1}\phi_y(f)\right|$$

$$\leq 2 \sup_{f \in \mathcal{F}} \max_{\lambda \in \Delta_m} \left|\sum_{y=1}^m \lambda_y \hat{\phi}_y(f) - \sum_{y=1}^m \lambda_y\phi_y(f)\right|$$

$$\leq 2 \sup_{f \in \mathcal{F}} \max_{\lambda \in \Delta_m} \sum_{y=1}^m \lambda_y \left|\hat{\phi}_y(f) - \phi_y(f)\right|$$

$$= 2 \sup_{f \in \mathcal{F}} \max_{y \in [m]} \left|\hat{\phi}_y(f) - \phi_y(f)\right|.$$

where for the second inequality, we use (5). Applying the generalization bound assumed in Theorem 4, we have with probability $\leq 1 - \delta$ (over draw of $S \sim D^n$), for all iterations $k \in [K]$,

$$\sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(\tilde{f}) \leq 2\Delta(n, \delta),$$

For the second part, note that for any $\lambda \in \Delta_m$,

$$
\left| \sum_{y=1}^{m} \lambda_y \hat{\phi}_y^{\mathrm{val}}(f^{k+1}) - \sum_{y=1}^{m} \lambda_y \phi_y(f^{k+1}) \right| \leq \sum_{y=1}^{m} \lambda_y \left| \hat{\phi}_y^{\mathrm{val}}(f^{k+1}) - \phi_y(f^{k+1}) \right|
$$

$$
\leq \max_{y \in [m]} \left| \hat{\phi}_y^{\mathrm{val}}(f^{k+1}) - \phi_y(f^{k+1}) \right|
$$

$$
\leq \sup_{f \in \mathcal{F}} \max_{y \in [m]} \left| \hat{\phi}_y^{\mathrm{val}}(f) - \phi_y(f) \right|.
$$

An application of the generalization bound assumed in Theorem 4 to empirical estimates from the validation sample completes the proof. □

We are now ready to prove Theorem 4.

*Proof of Theorem 4.* Note that because $\min_{y \in [m]} \pi_y^t \geq \frac{1}{Z}$ and $n^{\mathrm{val}} \geq 8Z \log(2m/\delta)$, we have by a direct application of Chernoff's bound (along with a union bound over all $m$ classes) that with probability at least $1 - \delta/2$:

$$\min_{y \in [m]} \hat{\pi}_y^{t,\mathrm{val}} \geq \frac{1}{2Z}, \forall y \in [m]$$

and consequently, $\hat{\phi}_y^{\mathrm{val}}(f) \leq 2BZ, \forall f \in \mathcal{F}$. The boundedness of $\hat{\phi}_y^{\mathrm{val}}$ will then allow us to apply standard convergence guarantees for exponentiated gradient ascent [Shalev-Shwartz et al., 2011]. For $\gamma = \frac{1}{2BZ}\sqrt{\frac{\log(m)}{K}}$, the updates on $\lambda$ will give us with probability at least $1 - \delta/2$:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y \hat{\phi}_y^{\mathrm{val}}(f^k) \leq \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y^k \hat{\phi}_y^{\mathrm{val}}(f^k) + 4BZ\sqrt{\frac{\log(m)}{K}} \tag{6}$$

Applying the second part of Lemma 2 to each iteration $k$, we have with probability at least $1 - \delta$:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y \phi_y(f^k) \leq \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y^k \phi_y(f^k) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\mathrm{val}}, \delta/2),$$

and applying the first part of Lemma 2 to the RHS, we have with the same probability:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y \phi_y(f^k)$$

$$\leq \frac{1}{K} \sum_{k=1}^{K} \min_{f \in \mathcal{F}} \sum_{y=1}^{m} \lambda_y^k \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\mathrm{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$\leq \min_{f \in \mathcal{F}} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y^k \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\mathrm{val}}, \delta/2) + 2\Delta(n, \delta/2).$$

Note that we have taken a union bound over the high probability statement in (6) and that in Lemma 2. Using the convexity of $\phi(\cdot)$ in $f(x)$ and Jensen's inequality, we have that $\sum_{y=1}^{m} \lambda_y \phi_y(\bar{f}^s) \leq \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y \phi_y(f^k)$. We use this to further lower bound the LHS in terms of the averaged scoring function $\bar{f}^s(x) = \frac{1}{K} \sum_{k=1}^{K} f^k(x)$:

$$\max_{\lambda \in \Delta_m} \sum_{y=1}^{m} \lambda_y \phi_y(\bar{f}^s)$$

$$\leq \min_{f \in \mathcal{F}} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y^k \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$= \min_{f \in \mathcal{F}} \sum_{y=1}^{m} \tilde{\lambda}_y \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$\leq \max_{\lambda \in \Delta_m} \min_{f \in \mathcal{F}} \sum_{y=1}^{m} \lambda_y \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$= \min_{f \in \mathcal{F}} \max_{\lambda \in \Delta_m} \sum_{y=1}^{m} \lambda_y \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$= \min_{f \in \mathcal{F}} \max_{y \in [m]} \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2), \tag{7}$$

where in the second step $\tilde{\lambda}_y = \frac{1}{K}\sum_{k=1}^{K} \lambda_y^k$; in the fourth step, we swap the 'min' and 'max' using Sion's minimax theorem [Sion, 1958]. We further have from (7),

$$\max_{y \in [m]} \phi_y(\bar{f}^s) \leq \min_{f \in \mathcal{F}} \max_{y \in [m]} \phi_y(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2).$$

In other words,

$$L^{\text{rob-d}}(\bar{f}^s) \leq \min_{f \in \mathcal{F}} L^{\text{rob-d}}(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2).$$

To complete the proof, we need to turn this into a guarantee on the original robust objective $L^{\text{rob}}$ in (3):

$$L^{\text{rob}}(\bar{f}^s)$$

$$\leq \min_{f \in \mathcal{F}} L^{\text{rob}}(f) + 2\max_{f \in \mathcal{F}} \left| L^{\text{rob}}(f) - L^{\text{rob-d}}(f) \right| + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2)$$

$$\leq \min_{f \in \mathcal{F}} L^{\text{rob}}(f) + 2B\max_{y \in [m]} \mathbb{E}_x\left[\left|\frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right] + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta(n^{\text{val}}, \delta/2) + 2\Delta(n, \delta/2),$$

where we have used the bound on the approximation error in the proof of Theorem 2. This completes the proof. $\square$

# B    STUDENT ESTIMATION ERROR

We now provide a bound on the estimation error in Theorem 4 using a generalization bound from Menon et al. [2021a].

**Lemma 3.** *Let $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$ be a given class of scoring functions. Let $\mathcal{V} \subseteq \mathbb{R}^{\mathcal{X}}$ denote the class of loss functions $v(x, y) = \ell(y, f(x))$ induced by scorers $f \in \mathcal{F}$. Let $\mathcal{M}_n = \mathcal{N}_\infty(\frac{1}{n}, \mathcal{V}, 2n)$ denote the uniform $L_\infty$ covering number for $\mathcal{V}$. Fix $\delta \in (0, 1)$. Suppose $\ell(y, z) \leq B$, $\pi_y^t \leq \frac{1}{Z}, \forall y \in [m]$, and the number of samples $n \geq 8Z\log(4m/\delta)$. Then with probability $\geq 1 - \delta$ over draw of $S \sim D^n$, for any $f \in \mathcal{F}$ and $y \in [m]$:*

$$\left|\phi_y(f) - \hat{\phi}_y(f)\right| \leq CZ\left(\sqrt{\mathbb{V}_{n,y}(f)\frac{\log(m\mathcal{M}_n/\delta)}{n}} + \frac{\log(m\mathcal{M}_n/\delta)}{n} + B\sqrt{\frac{\log(m/\delta)}{n}}\right),$$

*where $\mathbb{V}_{n,y}(f)$ denotes the empirical variance of the loss values $\{p_y^t(x_i) \cdot \ell(y, f(x_i))\}_{i=1}^n$ for class $y$, and $C > 0$ is a distribution-independent constant.*

Notice the dependence on the *variance* that the teacher's predictions induce on the loss. This suggests that the lower the variance in the teacher's predictions, the better is the student's generalization. Similar to Menon et al. [2021a], one can further show that when the teacher closely approximates the Bayes-probabilities $\eta(x)$, the distilled loss $p_y^t(x_i) \cdot \ell(y, f(x_i))$ has a lower empirical variance that the loss $\ell(y_i, f(x_i))$ computed from one-hot labels.

*Proof of Lemma 3.* We begin by defining the following intermediate term:

$$\tilde{\phi}_y(f) = \frac{1}{\pi_y^t} \frac{1}{n} \sum_{i=1}^{n} p_y^t(x_i)\, \ell\left(y, f(x_i)\right).$$

Then for any $y \in [m]$,

$$\left|\phi_y(f) - \hat{\phi}_y(f)\right| \leq \left|\phi_y(f) - \tilde{\phi}_y(f)\right| + \left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right|. \tag{8}$$

We next bound each of the terms in (8), starting with the first term:

$$\left|\phi_y(f) - \tilde{\phi}_y(f)\right| = \frac{1}{\pi_y^t}\left|\mathbb{E}_x\left[p_y^t(x)\, \ell\left(y, f(x)\right)\right] - \frac{1}{n}\sum_{i=1}^{n} p_y^t(x_i)\, \ell\left(y, f(x_i)\right)\right|$$

$$\leq Z\left|\mathbb{E}_x\left[p_y^t(x)\, \ell\left(y, f(x)\right)\right] - \frac{1}{n}\sum_{i=1}^{n} p_y^t(x_i)\, \ell\left(y, f(x_i)\right)\right|,$$

where we use the fact that $\pi_y^t \leq \frac{1}{Z}, \forall y$. Applying the generalization bound from Menon et al. [2021a, Proposition 2], along with a union bound over all $m$ classes, we have with probability at least $1 - \delta/2$ over the draw of $S \sim D^n$, for all $y \in [m]$:

$$\left|\phi_y(f) - \tilde{\phi}_y(f)\right| \leq C' Z\left(\sqrt{\mathbb{V}_{n,y}(f)\frac{\log(m\mathcal{M}_n/\delta)}{n}} + \frac{\log(m\mathcal{M}_n/\delta)}{n}\right), \tag{9}$$

for a distribution-independent constant $C' > 0$.

We next bound the second term in (8):

$$\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| = \left|\frac{1}{\pi_y^t} - \frac{1}{\hat{\pi}_y^t}\right|\frac{1}{n}\sum_{i=1}^{n} p_y^t(x_i) \cdot \ell\left(y, f(x_i)\right)$$

$$\leq B\left|\frac{1}{\pi_y^t} - \frac{1}{\hat{\pi}_y^t}\right|$$

$$= \frac{B}{\pi_y^t \hat{\pi}_y^t}\left|\pi_y^t - \hat{\pi}_y^t\right|,$$

where in the second step we use the fact that $\ell(y, f(x)) \leq B$ and $p_y^t(x) \leq 1$.

Further note that because $\min_{y \in [m]} \pi_y^t \geq \frac{1}{Z}$ and $n \geq 8Z\log(4m/\delta)$, we have by a direct application of Chernoff's bound (and a union bound over $m$ classes) that with probability at least $1 - \delta/4$:

$$\min_{y \in [m]} \hat{\pi}_y^t \geq \frac{1}{2Z}, \forall y \in [m]. \tag{10}$$

Therefore for any $y \in [m]$:

$$\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \leq 2BZ^2\left|\pi_y^t - \hat{\pi}_y^t\right|.$$

Conditioned on the above statement, a simple application of Hoeffding's inequality and a union bound over all $y \in [m]$ gives us that with probability at least $1 - \delta/4$ over the draw of $S \sim D^n$, for all $y \in [m]$:

$$\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \leq 2BZ^2\left(\frac{1}{Z}\sqrt{\frac{\log(8m/\delta)}{2n}}\right) = 2BZ\sqrt{\frac{\log(8m/\delta)}{2n}}. \tag{11}$$

A union bound over the high probability statements in (9–11) completes the proof. To see this, note that, for any $\epsilon > 0$ and $y \in [m]$,

$$\mathbb{P}\left(\left|\phi_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon\right)$$

**Algorithm 1** Distilled Margin-based DRO with One-hot Validation Labels

---

**Inputs:** Teacher $p^t$, Student hypothesis class $\mathcal{F}$, Training set $S$, Validation set $S^{\text{val}}$, Step-size $\gamma \in \mathbb{R}_+$, Number of iterations $K$, Loss $\ell$
**Initialize:** Student $f^0 \in \mathcal{F}$, Multipliers $\lambda^0 \in \Delta_m$
**For** $k = 0$ to $K - 1$

$\qquad \tilde{\lambda}_j^{k+1} = \lambda_j^k \exp\left(\gamma \hat{R}_j\right), \forall j \in [m]$

$\qquad\qquad$ where $\hat{R}_j = \dfrac{1}{n^{\text{val}}} \dfrac{1}{\hat{\pi}_j^{\text{val}}} \displaystyle\sum_{(x,y) \in S^{\text{val}}} \ell(y, f^k(x))$ and $\hat{\pi}_j^{\text{val}} = \dfrac{1}{n^{\text{val}}} \displaystyle\sum_{(x,y) \in S^{\text{val}}} \mathbf{1}(y = j)$

$\qquad \lambda_y^{k+1} = \dfrac{\tilde{\lambda}_y^{k+1}}{\sum_{j=1}^m \tilde{\lambda}_j^{k+1}}, \forall y$

$\qquad f^{k+1} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \dfrac{1}{n} \displaystyle\sum_{i=1}^n \mathcal{L}^{\text{mar}}\left(p^t(x_i), f(x_i); \dfrac{\lambda^{k+1}}{\hat{\pi}^t}\right)$    // Replaced with a few steps of SGD

**End For**
**Output:** $\bar{f}^s : x \mapsto \frac{1}{K} \sum_{k=1}^K f^k(x)$

---

$$
\begin{aligned}
&\leq \mathbb{P}\left(\left(\left|\phi_y(f) - \tilde{\phi}_y(f)\right| \geq \epsilon\right) \vee \left(\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon\right)\right) \\
&\leq \mathbb{P}\left(\left|\phi_y(f) - \tilde{\phi}_y(f)\right| \geq \epsilon\right) + \mathbb{P}\left(\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon\right) \\
&\leq \mathbb{P}\left(\left|\phi_y(f) - \tilde{\phi}_y(f)\right| \geq \epsilon\right) + \mathbb{P}\left(\hat{\pi}_y^t \leq \frac{1}{Z}\right) \cdot \mathbb{P}\left(\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon \,\middle|\, \hat{\pi}_y^t \leq \frac{1}{Z}\right) \\
&\qquad\qquad\qquad\qquad + \mathbb{P}\left(\hat{\pi}_y^t \geq \frac{1}{Z}\right) \cdot \mathbb{P}\left(\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon \,\middle|\, \hat{\pi}_y^t \geq \frac{1}{Z}\right) \\
&\leq \mathbb{P}\left(\left|\phi_y(f) - \tilde{\phi}_y(f)\right| \geq \epsilon\right) + \mathbb{P}\left(\hat{\pi}_y^t \leq \frac{1}{Z}\right) + \mathbb{P}\left(\left|\tilde{\phi}_y(f) - \hat{\phi}_y(f)\right| \geq \epsilon \,\middle|\, \hat{\pi}_y^t \geq \frac{1}{Z}\right),
\end{aligned}
$$

which implies that a union bound over (9–11) would give us the desired result in Lemma 3. $\qquad\square$

## C    DRO WITH ONE-HOT VALIDATION LABELS

The updates on $\lambda$ in Algorithm 1 use a validation set labeled by the teacher. One could instead perform these updates with a curated validation set containing the original one-hot labels. Each of these choices presents different merits. The use of a teacher-labeled validation set is useful in many real world scenarios where labeled data is hard to obtain, while unlabeled data abounds. In contrast, the use of one-hot validation labels, although more expensive to obtain, may make the student more immune to errors in the teacher's predictions, as the coefficients $\lambda$s are now based on an unbiased estimate of the student's performance on each class.

Algorithm 1 contains a version of the margin-based DRO described in Section 3.1, where instead of teacher labels the original one-hot labels are used in the validation set.

Before proceeding to providing a convergence guarantee for this algorithm, we will find it useful to define the following one-hot metrics:

$$
\phi_y^{\text{oh}}(f^s) = \frac{1}{\pi_y} \mathbb{E}_x \left[\eta_y(x) \, \ell\left(y, f^s(x)\right)\right]
$$

$$
\hat{\phi}_y^{\text{oh,val}}(f^s) = \frac{1}{\hat{\pi}_y} \frac{1}{n^{\text{val}}} \sum_{(x',y') \in S^{\text{val}}} \mathbf{1}(y' = y) \, \ell\left(y', f^s(x')\right).
$$

**Theorem 4.** *Suppose $p^t \in \mathcal{F}$ and $\mathcal{F}$ is closed under linear transformations. Then the final scoring function $\bar{f}^s(x) = \frac{1}{K} \sum_{k=1}^K f^k(x)$ output by Algorithm 1 is of the form:*

$$
\operatorname{softmax}_y(\bar{f}^s(x')) \propto \bar{\lambda}_y p_y^t(x'), \quad \forall (x', y') \in S,
$$

where $\bar{\lambda}_y = \left(\prod_{k=1}^K \lambda_y^k / \pi_y^t\right)^{1/K}$. *Furthermore, suppose $\ell$ is the softmax cross-entropy loss in $\ell^{\mathrm{xent}}$, $\ell(y, z) \leq B$, for some $B > 0$, and $\max_{y \in [m]} \frac{1}{\pi_y} \leq Z$, for some $Z > 0$. Suppose for any $\delta \in (0, 1)$, the following holds: with probability at least $1 - \delta$ (over draw of $S \sim D^n$), for all $f \in \mathcal{F}$,*

$$\max_{y \in [m]} \left|\phi_y^{\mathrm{oh}}(f) - \hat{\phi}_y^{\mathrm{oh}}(f)\right| \leq \Delta^{\mathrm{oh}}(n, \delta); \qquad \max_{y \in [m]} \left|\phi_y(f) - \hat{\phi}_y(f)\right| \leq \Delta(n, \delta),$$

*for some $\Delta^{\mathrm{oh}}(n, \delta), \Delta(n, \delta) \in \mathbb{R}_+$ that is increasing in $1/\delta$, and goes to 0 as $n \to \infty$. Fix $\delta \in (0, 1)$. Then when the step size $\gamma = \frac{1}{2BZ}\sqrt{\frac{\log(m)}{K}}$ and $n^{\mathrm{val}} \geq 8Z \log(2m/\delta)$, with probability at least $1 - \delta$ (over draw of $S \sim D^n$ and $S^{\mathrm{val}} \sim D^{n^{\mathrm{val}}}$), for any $\tau \in \mathbb{R}_+$,*

$$L^{\mathrm{rob}}(\bar{f}^s) \leq \min_{f \in \mathcal{F}} L^{\mathrm{rob}}(f) + \underbrace{2B \max_{y \in [m]} \mathbb{E}_x \left[\left|\tau \cdot \frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right]}_{\text{Approximation error}}$$

$$+ \underbrace{2\tau \cdot \Delta^{\mathrm{oh}}(n^{\mathrm{val}}, \delta/2) + 2\Delta(n, \delta/2)}_{\text{Estimation error}} + \underbrace{4BZ\sqrt{\frac{\log(m)}{K}}}_{\text{EG convergence}}.$$

Comparing this to the bound in Theorem 4, we can see that there is an additional scaling factor $\tau$ against the teacher probabilities $p_y^t(x)$ and in the approximation error. When we set $\tau = 1$, the bound looks very similar to Theorem 4, except that the estimation error term $\Delta^{\mathrm{oh}}$ now involves one-hot labels. Therefore the estimation error may incur a slower convergence with sample size as it no longer benefits from the lower variance that the teacher predictions may offer (see Appendix B for details).

The $\tau$-scaling in the approximation error also means that the teacher is no longer required to exactly match the (normalized) class probabilities $\eta(x)$. In fact, one can set $\tau$ to a value for which the approximation error is the lowest, and in general to a value that minimizes the upper bound in Theorem 4, potentially providing us with a tighter convergence rate than Theorem 4.

The proof of Theorem 4 is similar to that of Theorem 4, but requires a modified version of Lemma 2:

**Lemma 5.** *Suppose the conditions in Theorem 4 hold. With probability $\leq 1 - \delta$ (over draw of $S \sim D^n$ and $S^{\mathrm{val}} \sim D^{n^{\mathrm{val}}}$), at each iteration $k$ and for any $\tau \in \mathbb{R}_+$,*

$$\sum_{y=1}^m \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(f^{k+1}) - \min_{f \in \mathcal{F}} \sum_{y=1}^m \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(f) \leq 2\tau \cdot \Delta(n, \delta) + 2B \max_{y \in [m]} \mathbb{E}_x \left[\left|\tau \frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right].$$

*Furthermore, with the same probability, for any $\lambda \in \Delta_m$:*

$$\left|\sum_{y=1}^m \lambda_y \hat{\phi}_y^{\mathrm{oh,val}}(f^{k+1}) - \sum_{y=1}^m \lambda_y \phi_y^{\mathrm{oh}}(f^{k+1})\right| \leq \Delta^{\mathrm{oh}}(n^{\mathrm{val}}, \delta).$$

*Proof.* We first note from Lemma 1 that because $f^{k+1} \in \operatorname*{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{\mathrm{mar}}\left(p^t(x_i), f(x_i); \frac{\lambda^{k+1}}{\hat{\pi}}\right)$, we have for the example-weighting $w_i = \tau, \forall i$:

$$\tau \sum_{y=1}^m \lambda_y^{k+1} \hat{\phi}_y(f^{k+1}) \leq \tau \sum_{y=1}^m \lambda_y^{k+1} \hat{\phi}_y(f), \ \forall f \in \mathcal{F}. \tag{12}$$

For a fixed iteration $k$, let us denote $\tilde{f} \in \operatorname*{argmin}_{f \in \mathcal{F}} \sum_{y=1}^m \lambda_y^{k+1} \phi_y(f)$. Then for the first part, we have for any $\tau \in \mathbb{R}_+$:

$$\sum_{y=1}^m \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(f^{k+1}) - \sum_{y=1}^m \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(\tilde{f})$$

$$\leq \tau \left( \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(\tilde{f}) \right) + \sum_{y=1}^{m} \lambda_y^{k+1} \left| \phi_y^{\mathrm{oh}}(f^{k+1}) - \tau \phi_y(f^{k+1}) \right|$$

$$+ \sum_{y=1}^{m} \lambda_y^{k+1} \left| \phi_y^{\mathrm{oh}}(\tilde{f}) - \tau \phi_y(\tilde{f}) \right|$$

$$\leq \tau \left( \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(\tilde{f}) \right) + 2 \max_{f \in \mathcal{F}} \sum_{y=1}^{m} \lambda_y^{k+1} \left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right|$$

$$\leq \tau \left( \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(\tilde{f}) \right) + 2 \max_{f \in \mathcal{F}} \max_{\lambda \in \Delta_m} \sum_{y=1}^{m} \lambda \left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right|$$

$$\leq \tau \left( \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y(\tilde{f}) \right) + 2 \max_{f \in \mathcal{F}} \max_{y \in [m]} \left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right|$$

$$\leq 2\tau \sup_{f \in \mathcal{F}} \max_{y \in [m]} \left| \hat{\phi}_y(f) - \phi_y(f) \right| + 2 \max_{f \in \mathcal{F}} \max_{y \in [m]} \left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right|.$$

where the last inequality re-traces the steps in Lemma 2. Further applying the generalization bound assumed in Theorem 4, we have with probability $\leq 1 - \delta$ (over draw of $S \sim D^n$), for all iterations $k \in [K]$ and any $\tau \in \mathbb{R}_+$,

$$\sum_{y=1}^{m} \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(f^{k+1}) - \sum_{y=1}^{m} \lambda_y^{k+1} \phi_y^{\mathrm{oh}}(\tilde{f}) \leq 2\tau \Delta(n, \delta) + 2 \max_{f \in \mathcal{F}} \max_{y \in [m]} \left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right|. \tag{13}$$

All that remains is to bound the second term in (13). For any $f \in \mathcal{F}$ and $y \in [m]$,

$$\left| \phi_y^{\mathrm{oh}}(f) - \tau \phi_y(f) \right| \leq \left| \frac{1}{\pi_y} \mathbb{E}_x \left[ \eta_y(x) \, \ell \left( y, f(x) \right) \right] - \frac{\tau}{\pi_y^t} \mathbb{E}_x \left[ p_y^t(x) \, \ell \left( y, f(x) \right) \right] \right|$$

$$\leq \mathbb{E}_x \left[ \left| \frac{1}{\pi_y} \eta_y(x) \, \ell \left( y, f(x) \right) - \frac{\tau}{\pi_y^t} p_y^t(x) \, \ell \left( y, f(x) \right) \right| \right]$$

$$= \mathbb{E}_x \left[ \left| \frac{1}{\pi_y} \eta_y(x) - \frac{\tau}{\pi_y^t} p_y^t(x) \right| \ell \left( y, f^s(x) \right) \right]$$

$$\leq B \mathbb{E}_x \left[ \left| \frac{\eta_y(x)}{\pi_y} - \tau \frac{p_y^t(x)}{\pi_y^t} \right| \right],$$

where we use Jensen's inequality in the second step, the fact that $\ell(y, z) \leq B$ is non-negative in the second step, and the fact that $\ell(y, z) \leq B$ in the last step. Substituting this upper bound back into (13) completes the proof of the first part.

The second part follows from a direct application of the bound on the per-class estimation error $\max_{y \in [m]} \left| \phi_y^{\mathrm{oh}}(f) - \hat{\phi}_y^{\mathrm{oh,val}}(f) \right|$. $\qquad \square$

*Proof of Theorem 4.* The proof traces the same steps as Proposition 3 and Theorem 4, except that it applies Lemma 5 instead of Lemma 2.

Note that because $\min_{y \in [m]} \pi_y \geq \frac{1}{Z}$ and $n^{\mathrm{val}} \geq 8Z \log(2m/\delta)$, we have by a direct application of Chernoff's bound (along with a union bound over all $m$ classes) that with probability at least $1 - \delta/2$:

$$\min_{y \in [m]} \hat{\pi}_y^{\mathrm{oh,val}} \geq \frac{1}{2Z}, \forall y \in [m],$$

and consequently, $\hat{\phi}_y^{\mathrm{oh,val}}(f) \leq 2BZ, \forall f \in \mathcal{F}$. The boundedness of $\hat{\phi}_y^{\mathrm{oh,val}}$ will then allow us to apply standard convergence guarantees for exponentiated gradient ascent [Shalev-Shwartz et al., 2011]. For $\gamma = \frac{1}{2BZ} \sqrt{\frac{\log(m)}{K}}$, the updates on $\lambda$ will give us:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y \hat{\phi}_y^{\mathrm{oh,val}}(f^k) \leq \frac{1}{K} \sum_{k=1}^{K} \sum_{y=1}^{m} \lambda_y^k \hat{\phi}_y^{\mathrm{oh,val}}(f^k) + 4BZ \sqrt{\frac{\log(m)}{K}}$$

**Algorithm 2** Distilled Margin-based DRO for Traded-off Objective

---

**Inputs:** Teacher $p^t$, Student hypothesis class $\mathcal{F}$, Training set $S$, Validation set $S^{\text{val}}$, Step-size $\gamma \in \mathbb{R}_+$, Number of iterations $K$, Loss $\ell$, Trade-off parameter $\alpha$

**Initialize:** Student $f^0 \in \mathcal{F}$, Multipliers $\lambda^0 \in \Delta_m$

**For** $k = 0$ to $K - 1$

$$\tilde{\lambda}_j^{k+1} = \lambda_j^k \exp\left(\gamma \alpha \hat{R}_j\right), \forall j \in [m] \text{ where } \hat{R}_j = \frac{1}{n^{\text{val}}} \frac{1}{\hat{\pi}_j^{t,\text{val}}} \sum_{(x,y) \in S^{\text{val}}} p_j^t(x_i)\, \ell(j, f^k(x))$$

$$\lambda_y^{k+1} = \frac{\tilde{\lambda}_y^{k+1}}{\sum_{j=1}^m \tilde{\lambda}_j^{k+1}}, \forall y$$

$$\beta_y^{k+1} = (1 - \alpha)\frac{1}{m} + \alpha \lambda_y^{k+1}$$

$$f^{k+1} \in \operatorname*{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{\text{mar}}\left(p^t(x_i), f(x_i); \frac{\beta^{k+1}}{\hat{\pi}^t}\right) \quad \text{// Replaced with a few steps of SGD}$$

**End For**

**Output:** $\bar{f}^s : x \mapsto \frac{1}{K} \sum_{k=1}^K f^k(x)$

---

Applying the second part of Lemma 2 to each iteration $k$, we have with probability at least $1 - \delta$:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^K \sum_{y=1}^m \lambda_y \phi_y^{\text{oh}}(f^k) \leq \frac{1}{K} \sum_{k=1}^K \sum_{y=1}^m \lambda_y^k \phi_y^{\text{oh}}(f^k) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta^{\text{oh}}(n^{\text{val}}, \delta/2),$$

and applying the first part of Lemma 2 to the RHS, we have with the same probability, for any $\tau \in \mathbb{R}_+$:

$$\max_{\lambda \in \Delta_m} \frac{1}{K} \sum_{k=1}^K \sum_{y=1}^m \lambda_y \phi_y^{\text{oh}}(f^k) \leq \frac{1}{K} \sum_{k=1}^K \min_{f \in \mathcal{F}} \sum_{y=1}^m \lambda_y^k \phi_y^{\text{oh}}(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta^{\text{oh}}(n^{\text{val}}, \delta/2)$$

$$+ 2\tau\Delta(n, \delta/2) + 2B \max_{y \in [m]} \mathbb{E}_x\left[\left|\tau \frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right]$$

$$\leq \min_{f \in \mathcal{F}} \frac{1}{K} \sum_{k=1}^K \sum_{y=1}^m \lambda_y^k \phi_y^{\text{oh}}(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta^{\text{oh}}(n^{\text{val}}, \delta/2)$$

$$+ 2\tau\Delta(n, \delta/2) + 2B \max_{y \in [m]} \mathbb{E}_x\left[\left|\tau \frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right].$$

Using the convexity of $\phi(\cdot)$ in $f(x)$ and Jensen's inequality, we have that $\sum_{y=1}^m \lambda_y \phi_y(\bar{f}^s) \leq \frac{1}{K} \sum_{k=1}^K \sum_{y=1}^m \lambda_y \phi_y(f^k)$. We use this to further lower bound the LHS in terms of the averaged scoring function $\bar{f}^s(x) = \frac{1}{K} \sum_{k=1}^K f^k(x)$, and re-trace the steps in Theorem 4 to get"

$$\max_{y \in [m]} \phi_y^{\text{oh}}(\bar{f}^s) \leq \min_{f \in \mathcal{F}} \max_{y \in [m]} \phi_y^{\text{oh}}(f) + 4BZ\sqrt{\frac{\log(m)}{K}} + 2\Delta^{\text{oh}}(n^{\text{val}}, \delta/2)$$

$$+ 2\tau\Delta(n, \delta/2) + 2B \max_{y \in [m]} \mathbb{E}_x\left[\left|\tau \frac{p_y^t(x)}{\pi_y^t} - \frac{\eta_y(x)}{\pi_y}\right|\right].$$

Noting that $L^{\text{rob}}(f) = \max_{y \in [m]} \phi_y^{\text{oh}}(f)$ completes the proof. $\qquad\square$

## D   DRO FOR TRADED-OFF OBJECTIVE

We present a variant of the margin-based DRO algorithm described in Section 3.1 that seeks to minimize a trade-off between the balanced and robust student objectives:

$$\hat{L}^{\text{tdf-d}}(f^s) = (1 - \alpha)\hat{L}^{\text{bal-d}}(f^s) + \alpha \hat{L}^{\text{rob-d}}(f^s),$$

for some $\alpha \in [0,1]$.

Expanding this, we have:

$$L^{\text{tdf-d}}(f) = (1-\alpha)\frac{1}{m}\sum_{y=1}^{m}\frac{1}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i)) \;+\; \alpha \max_{y\in[m]}\sum_{y=1}^{m}\frac{1}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i))$$

$$= (1-\alpha)\frac{1}{m}\sum_{y=1}^{m}\frac{1}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i)) \;+\; \alpha \max_{\lambda\in\Delta_m}\sum_{y=1}^{m}\frac{\lambda_y}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i))$$

$$= \max_{\lambda\in\Delta_m}\sum_{y=1}^{m}\left((1-\alpha)\frac{1}{m} + \alpha\lambda_y\right)\frac{1}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i)).$$

The minimization of $L^{\text{tdf-d}}(f)$ over $f$ can then be a cast as a min-max problem:

$$\min_{f:\mathcal{X}\to\mathbb{R}^m} L^{\text{tdf-d}}(f) = \min_{f:\mathcal{X}\to\mathbb{R}^m}\max_{\lambda\in\Delta_m}\sum_{y=1}^{m}\left((1-\alpha)\frac{1}{m} + \alpha\lambda_y\right)\frac{1}{\hat{\pi}_y^t}\frac{1}{n}\sum_{i=1}^{n}p_y^t(x_i)\,\ell(y, f(x_i)).$$

Retracing the steps in the derivation of Algorithm 1 in Section 3.1, we have the following updates on $\lambda$ and $f$ to solve the above min-max problem:

$$\tilde{\lambda}_y^{k+1} = \lambda_y^k \exp\left(\gamma\alpha\frac{1}{n\hat{\pi}_y^t}\sum_{i=1}^{n}p_y^t(x_i)\,\ell\left(y, f^k(x_i)\right)\right), \forall y$$

$$\lambda_y^{k+1} = \frac{\tilde{\lambda}_y^{k+1}}{\sum_{j=1}^{m}\tilde{\lambda}_j^{k+1}}, \forall y$$

$$\beta_y^{k+1} = (1-\alpha)\frac{1}{m} + \alpha\lambda_y^{k+1}$$

$$f^{k+1} \in \operatorname*{argmin}_{f\in\mathcal{F}}\sum_{y\in[m]}\frac{\beta_y^{k+1}}{n\hat{\pi}_y^t}\sum_{i=1}^{n}p_y^t(x_i)\,\ell\left(y, f(x_i)\right),$$

for step-size parameter $\gamma > 0$. To better handle training of over-parameterized students, we will perform the updates on $\lambda$ using a held-out validation set, and employ a margin-based surrogate for performing the minimization over $f$. This procedure is outlined in Algorithm 2.

## D.1   CONNECTION TO POST-HOC ADJUSTMENT

The form of the student in Proposition 3 raises an interesting question. Instead of training an explicit student model, why not directly construct a new scoring model by making post-hoc adjustments to the teacher's predictions? Specifically, one could optimize over functions of the form $f_y^s(x) = \log(\gamma_y p_y^t(x))$, where the teacher $p^t$ is fixed, and pick the coefficients $\gamma \in \mathbb{R}^m$ so that resulting scoring function yields the best worst-class accuracy on a held-out dataset. This simple *post-hoc adjustment* strategy may not be feasible if the goal is to distill to a student that is considerably smaller than the teacher. Often, this is the case in settings where distillation is used as a compression technique. Yet, this post-hoc method serves as good baseline to compare with.

# E   ADDITIONAL EXPERIMENT DETAILS

This section contains further experiment details about the datasets, hyperparameters, and baselines.

## E.1   ADDITIONAL DETAILS ABOUT DATASETS

### E.1.1   Building long tailed datasets

The long-tailed datasets were created from the original datasets following Cui et al. [2019] by downsampling examples with an exponential decay in the per-class sizes. As done by Narasimhan and Menon [2021], we set the imbalance ratio

$\frac{\max_i P(y=i)}{\min_i P(y=i)}$ to 100 for CIFAR-10 and CIFAR-100, and to 83 for TinyImageNet (the slightly smaller ratio is to ensure that the smallest class is of a reasonable size). We use the long-tail version of ImageNet generated by Liu et al. [2017].

### E.1.2 Dataset splits

The original test samples for CIFAR-10, CIFAR-10-LT, CIFAR-100, CIFAR-100-LT, TinyImageNet (200 classes), TinyImageNet-LT (200 classes), and ImageNet (1000 classes) are all balanced. Following Narasimhan and Menon [2021], we randomly split them in half and use half the samples as a validation set, and the other half as a test set. For the CIFAR and TineImageNet datasets, this amounts to using a validation set of size 5000. For the ImageNet dataset, we sample a subset of 5000 examples from the validation set each time we update the Lagrange multipliers in Algorithm 1.

In keeping with prior work Menon et al. [2021b], Narasimhan and Menon [2021], Lukasik et al. [2022], we use the same validation and test sets for the long-tailed training sets as we do for the original versions. For the long tailed training sets, this simulates a scenario where the training data follows a long tailed distribution due to practical data collection limitations, but the test distribution of interest still comes from the original data distribution. In plots, the "balanced accuracy" that we report for the long-tail datasets (e.g., CIFAR-10-LT) is actually the standard accuracy calculated over the balanced test set, which is shared with the original balanced dataset (e.g., CIFAR-10).

Both teacher and student were always trained on the same training set.

The CIFAR datasets had images of size $32 \times 32$, while the TinyImageNet and ImageNet datasets dataset had images of size $224 \times 224$.

These datasets do not contain personally identifiable information or offensive content. The CIFAR-10 and CIFAR-100 datasets are licensed under the MIT License. The terms of access for ImageNet are given at `https://www.image-net.org/download.php`.

### E.2 ADDITIONAL DETAILS ABOUT TRAINING AND HYPERPARAMETERS

### E.2.1 Training details and hyperparameters

**Temperature hyperparameters.** We apply temperature scaling to the teacher scores on both the training set and validation set when training the student, i.e., compute $p^t(x) = \text{softmax}(f^t(x)/\gamma)$, and vary the temperature parameter $\gamma$ over a range of $\{1, 3, 5\}$. When training with teacher labels on the validation set (Algorithm 1), we vary the temperature parameters independently for the training set and the validation set. That is, we apply $p^t(x) = \text{softmax}(f^t(x)/\gamma_{\text{train}})$ over the training set and $p^t(x) = \text{softmax}(f^t(x)/\gamma_{\text{val}})$ over the validation set. When teacher labels are applied to the validation set, we additionally include a temperature of 0.1 on the teacher's validation set labels to approximate a hard thresholding of the teacher probabilities. Thus, the final hyperparameter search spaces are $\gamma_{\text{train}} \in \{1, 3, 5\}$, and $\gamma_{\text{val}} \in \{0.1, 1, 3, 5\}$.

Unless otherwise specified, in all tables, the temperature hyperparameters were chosen to achieve the best worst-class accuracy on the validation set. In all scatter plots such as Figure 1, for each $\alpha^t, \alpha^s$ combination, temperature hyperparameters were selected to achieve the best worst-class accuracy on the validation set.

**Learning rate hyperparameters.** All models were trained using SGD with momentum of 0.9 [Lukasik et al., 2022, Narasimhan and Menon, 2021].

The learning rate schedule were chosen to mimic the settings in prior work Narasimhan and Menon [2021], Lukasik et al. [2022]. For CIFAR-10 and CIFAR-100 datasets, we ran the optimizer for 450 epochs, linearly warming up the learning rate till the 15th epoch, and then applied a step-size decay of 0.1 after the 200th, 300th and 400th epochs, as done by Lukasik et al. [2022]. For the long-tail versions of these datasets, we trained for 256 epochs, linearly warming up the learning rate till the 15th epoch, and then applied a step-size decay of 0.1 after the 96th, 192nd and 224th epochs, as done by Narasimhan and Menon [2021]. Similarly, for the TinyImageNet datasets, we train for 200 epochs, linearly warming up the learning rate till the 5th epoch, and then applying a decay of 0.1 after the 75th and 135th epochs, as done by Narasimhan and Menon [2021]. For ImageNet, we train for 90 epochs, linearly warming up the learning rate till the 5th epoch, then applying a decay of 0.1 after the 30th, 60th and 80th epochs, as done by Lukasik et al. [2022]. We used a batch size of 128 for the CIFAR-10 and the long-tailed TinyImageNet datasets [Narasimhan and Menon, 2021], a batch size of 512 for the balanced ImageNet dataset, a batch size of 2048 for the balanced TinyImageNet dataset, and a batch size of 1024 for other datasets Lukasik et al. [2022].

We apply an $L_2$ weight decay of $10^{-4}$ in all our SGD updates Lukasik et al. [2022]. This amounts to applying an $L_2$ *regularization* on the model parameters, and has the effect of keeping the model parameters (and as a result the loss function) bounded.

When training with the margin-based robust objective (see Algorithm 1), a separate step size $\alpha$ was applied for training the main model function $f$, and for updating the multipliers $\lambda$. We set $\alpha$ to 0.1 in all experiments.

**Hardware.** Model training was done using TPUv2.

### E.2.2 Repeats

For all comparative baselines without distillation (Group DRO, Post shift, and all teachers alone), we provide average results over $m$ retrained models ($m = 5$ for ImageNet / TinyImageNet, or $m = 10$ for CIFAR datasets). For students on all CIFAR* datasets, unless otherwise specified, we train the teacher once and run the student training 10 times using the same arbitrarily chosen fixed teacher. We compute the mean and standard error of metrics over these $m = 10$ runs. For the resource-heavy TinyImageNet and ImageNet students, we reduce the number of repeats to $m = 5$. This methodology captures variation in the student retrainings while holding the teacher fixed. To capture the end-to-end variation in both teacher and student training, we include Appendix F.4 and Table 3 which contains a rerun of the CIFAR experiments in Table 1 using a distinct teacher for each student retraining. The overall best teacher/student objective combinations did not change for most datasets, with the only exception coming from a difference in the use of validation set labels.

## E.3 ADDITIONAL DETAILS ABOUT ALGORITHMS AND BASELINES

### E.3.1 Practical improvements to Algorithms 1–2

Algorithms 1–2 currently return a scorer that averages over all $K$ iterates $\bar{f}^s(x) = \frac{1}{K} \sum_{k=1}^{K} f^k(x)$. While this averaging was required for our theoretical robustness guarantees to hold, in our experiments, we find it sufficient to simply return the last model $f^K$. Another practical improvement that we make to these algorithms following Cotter et al. [2019], is to employ the 0-1 loss while performing updates on $\lambda$, i.e., set $\ell = \ell^{0\text{-}1}$ in the $\lambda$-update step. We are able to do this because the convergence of the exponentiated gradient updates on $\lambda$ does not depend on $\ell$ being differentiable. This modification allows $\lambda$s to better reflect the model's per-class performance on the validation sample.

### E.3.2 Discussion on post-shifting baseline

We implement the post-shifting method in Narasimhan and Menon [2021] (Algorithm 3 in their paper), which provides for an efficient way to construct a scoring function of the form $f_y^s(x) = \log(\gamma_y p_y^t(x))$, for a fixed teacher $p^t$, where the coefficients $\gamma \in \mathbb{R}^m$ are chosen to maximize the worst-class accuracy on the validation dataset. Interestingly, in our experiments, we find this approach to do exceedingly well on the validation sample, but this does not always translate to good worst-class test performance. In contrast, some of the teacher-student combinations that we experiment with were seen to over-fit less to the validation sample, and as a result were able to generalize better to the test set. This could perhaps indicate that the teacher labels we use in these combinations benefit the student in a way that it improves its generalization. The variance reduction effect that Menon et al. [2021a] postulate may be one possible explanation for why we see this behavior.

## F ADDITIONAL EXPERIMENTAL RESULTS

This section contains additional experimental results.

### F.1 EXTENDED TABLES FOR OBJECTIVE COMBINATIONS

We include extended tables comparing worst-class performance for different combinations of teacher and student objectives. The mean and standard errors are reported over repeat trainings as described in Appendix E.2.2.

Table 1 is an extended version of Table 1 that includes standard errors for both worst-$k$ accuracy and average accuracy.

Table 2 includes similar comparisons when the student is compressed – that is, the student's architecture is smaller than the teacher's architecture.

Table 1: Worst-class accuracy comparison of self-distilled teacher/student combos on test. The "none" row indicates the performance of the teacher alone. Worst-class accuracy is shown above (or worst-10 accuracy for TinyImageNet-LT), and average is accuracy shown in parentheses below. The combination with the best worst-class accuracy is in **bold**. We include results for the robust student using either a teacher labeled validation set ("teacher val"), or true one-hot class labels in the validation set ("one-hot val"), as outlined in Appendix C. Perhaps counterintuitively, the teacher with the best worst-class accuracy alone (the "none" row) did not always produce the student with the highest worst-class accuracy.

| Student Obj. | **CIFAR-10** Teacher Obj. $L^{std}$ | $L^{rob}$ | **CIFAR-100** Teacher Obj. $L^{std}$ | $L^{rob}$ | **TinyImageNet** Teacher Obj. $L^{std}$ | $L^{rob}$ |
|---|---|---|---|---|---|---|
| none | $86.48 \pm 0.32$ | $90.09 \pm 0.22$ | $42.22 \pm 0.90$ | $43.42 \pm 1.03$ | $8.42 \pm 1.88$ | $11.87 \pm 1.74$ |
| | $(93.74 \pm 0.05)$ | $(92.67 \pm 0.09)$ | $72.42 \pm 0.16$ | $68.81 \pm 0.11$ | $(56.79 \pm 0.33)$ | $(48.40 \pm 0.15)$ |
| $L^{std-d}$ | $87.66 \pm 0.40$ | $90.12 \pm 0.23$ | $43.81 \pm 0.58$ | $\mathbf{48.20 \pm 1.15}$ | $6.32 \pm 2.31$ | $10.53 \pm 1.49$ |
| | $(94.34 \pm 0.07)$ | $(94.07 \pm 0.07)$ | $(74.61 \pm 0.15)$ | $(73.23 \pm 0.07)$ | $(57.83 \pm 0.13)$ | $(55.36 \pm 0.16)$ |
| $L^{rob-d}$ (teacher val) | $\mathbf{90.94 \pm 0.16}$ | $85.14 \pm 0.47$ | $39.18 \pm 1.58$ | $30.42 \pm 1.30$ | $9.98 \pm 1.87$ | $16.58 \pm 1.23$ |
| | $(92.54 \pm 0.05)$ | $(89.58 \pm 0.11)$ | $(63.49 \pm 0.29)$ | $(55.77 \pm 0.39)$ | $(49.84 \pm 0.21)$ | $(46.11 \pm 0.37)$ |
| $L^{rob-d}$ (one-hot val) | $89.37 \pm 0.17$ | $87.32 \pm 0.21$ | $44.61 \pm 1.55$ | $42.68 \pm 0.74$ | $16.27 \pm 0.43$ | $\mathbf{17.36 \pm 1.32}$ |
| | $(91.63 \pm 0.06)$ | $(91.16 \pm 0.10)$ | $(69.02 \pm 0.30)$ | $(62.03 \pm 0.24)$ | $(48.06 \pm 0.24)$ | $(43.92 \pm 0.30)$ |

| Student Obj. | **CIFAR-10-LT** Teacher Obj. $L^{std}$ | $L^{bal}$ | $L^{rob}$ | **CIFAR-100-LT** Teacher Obj. $L^{std}$ | $L^{bal}$ | $L^{rob}$ |
|---|---|---|---|---|---|---|
| None | $57.26 \pm 0.55$ | $68.52 \pm 0.52$ | $74.8 \pm 0.30$ | $0.00 \pm 0.00$ | $3.75 \pm 0.62$ | $10.33 \pm 0.82$ |
| | $(76.27 \pm 0.20)$ | $(79.85 \pm 0.20)$ | $(80.29 \pm 0.12)$ | $(43.33 \pm 0.16)$ | $(47.55 \pm 0.17)$ | $(44.27 \pm 0.13)$ |
| $L^{std-d}$ | $36.67 \pm 0.28$ | $66.96 \pm 0.43$ | $71.15 \pm 0.24$ | $0.00 \pm 0.00$ | $2.39 \pm 0.24$ | $7.32 \pm 0.47$ |
| | $(69.5 \pm 0.13)$ | $(79.25 \pm 0.10)$ | $(80.95 \pm 0.11)$ | $(43.86 \pm 0.14)$ | $(48.95 \pm 0.15)$ | $(47.93 \pm 0.11)$ |
| $L^{bal-d}$ | $71.23 \pm 0.44$ | $70.52 \pm 0.20$ | $72.96 \pm 0.53$ | $4.39 \pm 0.65$ | $7.08 \pm 0.80$ | $7.19 \pm 0.79$ |
| | $(80.5 \pm 0.12)$ | $(81.12 \pm 0.08)$ | $(80.71 \pm 0.07)$ | $(50.4 \pm 0.11)$ | $(50.1 \pm 0.09)$ | $(47.51 \pm 0.20)$ |
| $L^{rob-d}$ (teacher val) | $63.85 \pm 0.21$ | $\mathbf{75.56 \pm 0.19}$ | $69.21 \pm 0.45$ | $9.05 \pm 0.71$ | $12.52 \pm 0.98$ | $10.32 \pm 0.76$ |
| | $(76.81 \pm 0.08)$ | $(80.81 \pm 0.08)$ | $(76.72 \pm 0.19)$ | $(33.75 \pm 0.10)$ | $(34.05 \pm 0.09)$ | $(36.83 \pm 0.15)$ |
| $L^{rob-d}$ (one-hot val) | $73.59 \pm 0.25$ | $75.43 \pm 0.38$ | $74.7 \pm_{0.19}$ | $12.28 \pm 0.46$ | $11.94 \pm 0.80$ | $\mathbf{13.18 \pm 0.61}$ |
| | $(77.92 \pm 0.05)$ | $(79.02 \pm 0.07)$ | $(77.99 \pm 0.10)$ | $(30.79 \pm 0.18)$ | $(29.8 \pm 0.20)$ | $(31.88 \pm 0.20)$ |

| Student Obj. | **TinyImageNet-LT** Teacher Obj. $L^{std}$ | $L^{bal}$ | $L^{rob}$ |
|---|---|---|---|
| None | $0.00 \pm 0.00$ | $2.11 \pm 0.37$ | $4.92 \pm 0.66$ |
| | $(33.15 \pm 0.17)$ | $(35.96 \pm 0.12)$ | $(27.23 \pm 0.15)$ |
| $L^{std-d}$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ | $1.87 \pm 0.23$ |
| | $(26.05 \pm 0.18)$ | $(27.21 \pm 0.15)$ | $(25.34 \pm 0.13)$ |
| $L^{bal-d}$ | $0.20 \pm 0.18$ | $2.82 \pm 0.14$ | $4.77 \pm 0.41$ |
| | $(30.43 \pm 0.06)$ | $(39.41 \pm 0.15)$ | $(38.41 \pm 0.15)$ |
| $L^{rob-d}$ (teacher val) | $0.00 \pm 0.00$ | $4.93 \pm 0.38$ | $3.32 \pm 0.43$ |
| | $(22.66 \pm 0.08)$ | $(35.43 \pm 0.18)$ | $(25.11 \pm 0.17)$ |
| $L^{rob-d}$ (one-hot val) | $1.55 \pm 0.37$ | $6.11 \pm 0.39$ | $\mathbf{6.19 \pm 0.25}$ |
| | $(21.59 \pm 0.19)$ | $(28.24 \pm 0.17)$ | $(25.30 \pm 0.18)$ |

### F.1.1 Robust distillation with a onehot-labeled validation set

Tables 1 and 2 also include results when the robust student is trained using a validation set using onehot labels, as described in Appendix C. We report the accuracies for this robust student for different teachers trained with the standard, balanced, and robust objectives in the last rows of Tables 1 and 2 ($L^{rob-d}$ (one-hot val)). We compare these to the robust student trained using teacher labels on the validation set ($L^{rob-d}$ (teacher val)), which require less labeled data.

Table 2: Comparison of ResNet-56→ResNet-32 distilled teacher/student combos on test on CIFAR datasets. Worst-class accuracy shown above, and average accuracy shown in parentheses below. The combination with the best worst-class accuracy is bolded. Mean and standard error are reported over 10 repeats. We include results for the robust student using either a teacher labeled validation set ("teacher val"), or true one-hot class labels in the validation set ("one-hot val"), as outlined in Section 3.1.

| | | **CIFAR-10** Teacher Obj. | | **CIFAR-100** Teacher Obj. | |
| --- | --- | --- | --- | --- | --- |
| | | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ |
| Student Obj. (ResNet-32) | $L^{\text{std-d}}$ | $86.4 \pm 0.27$ | $89.56 \pm 0.20$ | $41.82 \pm 1.12$ | $\mathbf{45.7 \pm 1.13}$ |
| | | $(93.73 \pm 0.05)$ | $(93.38 \pm 0.05)$ | $(73.19 \pm 0.10)$ | $(71.42 \pm 0.22)$ |
| | $L^{\text{rob-d}}$ (teacher val) | $\mathbf{89.61 \pm 0.27}$ | $83.8 \pm 0.95$ | $38.94 \pm 2.61$ | $19.15 \pm 0.00$ |
| | | $(92.20 \pm 0.08)$ | $(88.71 \pm 0.24)$ | $(62.28 \pm 0.40)$ | $(52.9 \pm 0.00)$ |
| | $L^{\text{rob-d}}$ (one-hot val) | $87.92 \pm 0.23$ | $86.57 \pm 0.24$ | $33.19 \pm 1.29$ | $41.23 \pm 0.84$ |
| | | $(90.89 \pm 0.12)$ | $(90.54 \pm 0.11)$ | $(57.43 \pm 0.29)$ | $(61.14 \pm 0.24)$ |

| | | **CIFAR-10-LT** Teacher Obj. | | | **CIFAR-100-LT** Teacher Obj. | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ |
| Student Obj. (ResNet-32) | $L^{\text{std-d}}$ | $57.23 \pm 0.53$ | $66.80 \pm 0.25$ | $72.36 \pm 0.39$ | $0.00 \pm 0.00$ | $1.38 \pm 0.39$ | $7.99 \pm 0.48$ |
| | | $(75.76 \pm 0.12)$ | $(78.99 \pm 0.06)$ | $(80.74 \pm 0.09)$ | $(44.33 \pm 0.11)$ | $(47.28 \pm 0.13)$ | $(47.34 \pm 0.08)$ |
| | $L^{\text{bal-d}}$ | $71.37 \pm 0.50$ | $71.00 \pm 0.45$ | $72.17 \pm 0.40$ | $3.57 \pm 0.58$ | $4.28 \pm 0.45$ | $5.58 \pm 0.53$ |
| | | $(81.13 \pm 0.12)$ | $(81.12 \pm 0.15)$ | $(79.91 \pm 0.08)$ | $(49.21 \pm 0.10)$ | $(46.56 \pm 0.13)$ | $(48.58 \pm 0.09)$ |
| | $L^{\text{rob-d}}$ (teacher val) | $64.1 \pm 0.36$ | $73.51 \pm 0.33$ | $69.90 \pm 0.42$ | $10.24 \pm 0.71$ | $\mathbf{13.41 \pm 0.72}$ | $11.27 \pm 0.61$ |
| | | $(76.34 \pm 0.12)$ | $(80.10 \pm 0.10)$ | $(76.37 \pm 0.14)$ | $(33.55 \pm 0.16)$ | $(33.37 \pm 0.17)$ | $(36.14 \pm 0.19)$ |
| | $L^{\text{rob-d}}$ (one-hot val) | $72.65 \pm 0.27$ | $74.39 \pm 0.34$ | $\mathbf{74.45 \pm 0.26}$ | $10.93 \pm 0.65$ | $12.2 \pm 0.65$ | $12.93 \pm 0.62$ |
| | | $(77.69 \pm 0.11)$ | $(78.68 \pm 0.16)$ | $(77.97 \pm 0.10)$ | $(29.48 \pm 0.22)$ | $(30.27 \pm 0.18)$ | $(31.83 \pm 0.17)$ |

Perhaps surprisingly, it did not always benefit the robust student to utilize the true one-hot labels in the validation set. Instead, training the robust student with teacher labels on the validation set was often sufficient to achieve the best or close to the best worst-class performance. This is promising from a data efficiency standpoint, since it can be expensive to build up a labeled dataset for validation, especially if the training data is long-tailed.

## F.2 ADDITIONAL PLOTS FOR ALL TRADE-OFF PARAMETER COMBINATIONS

Figure 1 show accuracies for all $\alpha^t, \alpha^s$ the equivalent of Figure 1 but for all datasets.

## F.3 COMPARISON TO BASELINES OF ALL PARETO EFFICIENT TRADE-OFF PARAMETERS

To supplement the comparison to baselines in Table 2, Figures 2 and 3 show all Pareto efficient $\alpha^t$ and $\alpha^s$ combinations on test. Whereas only a single $\alpha^t, \alpha^s$ combination was selected on the validation set and reported in Table 2, Figures 2 and 3 show that there were many more combinations of $\alpha^t, \alpha^s$ that could have Pareto dominated all baselines.

Figures 2 and 3 also give more insight into which values of $\alpha^t$ work best for different values of $\alpha^s$. Whereas Figure 1 shows that $\alpha^s$ is highly correlated with average accuracy, the same is not true for $\alpha^t$. Worst-class accuracy generally increases with $\alpha^s$, but the teachers that achieve the Pareto efficient points all have $\alpha^t < 1$. This reveals counter-intuitively that the teacher's worst-class accuracy is not a direct predictor of the robustness of a subsequent student. This couples with our theoretical understanding in Section 5, which showed that the ability of a teacher to train robust students is determined by the calibration of scores within each class.

*Trading off average vs. worst-class accuracy.* Figures 2 and 3 show that when we allow for more nuanced $L^{\text{tdf}}$ objective combinations, the resulting models may have higher average accuracy and worst-class accuracy than standard distillation. Interestingly, the models with the most "even" trade-offs between average accuracy and worst-class accuracy tend to have low $\alpha^t$ (around 0.25) and low $\alpha^s$ (also around 0.25). Higher values of $\alpha^t$ tended to lead to more extreme points on the trade-off curve, either with higher average accuracy at the expense of worst-class accuracy, or vice versa. Overall, the
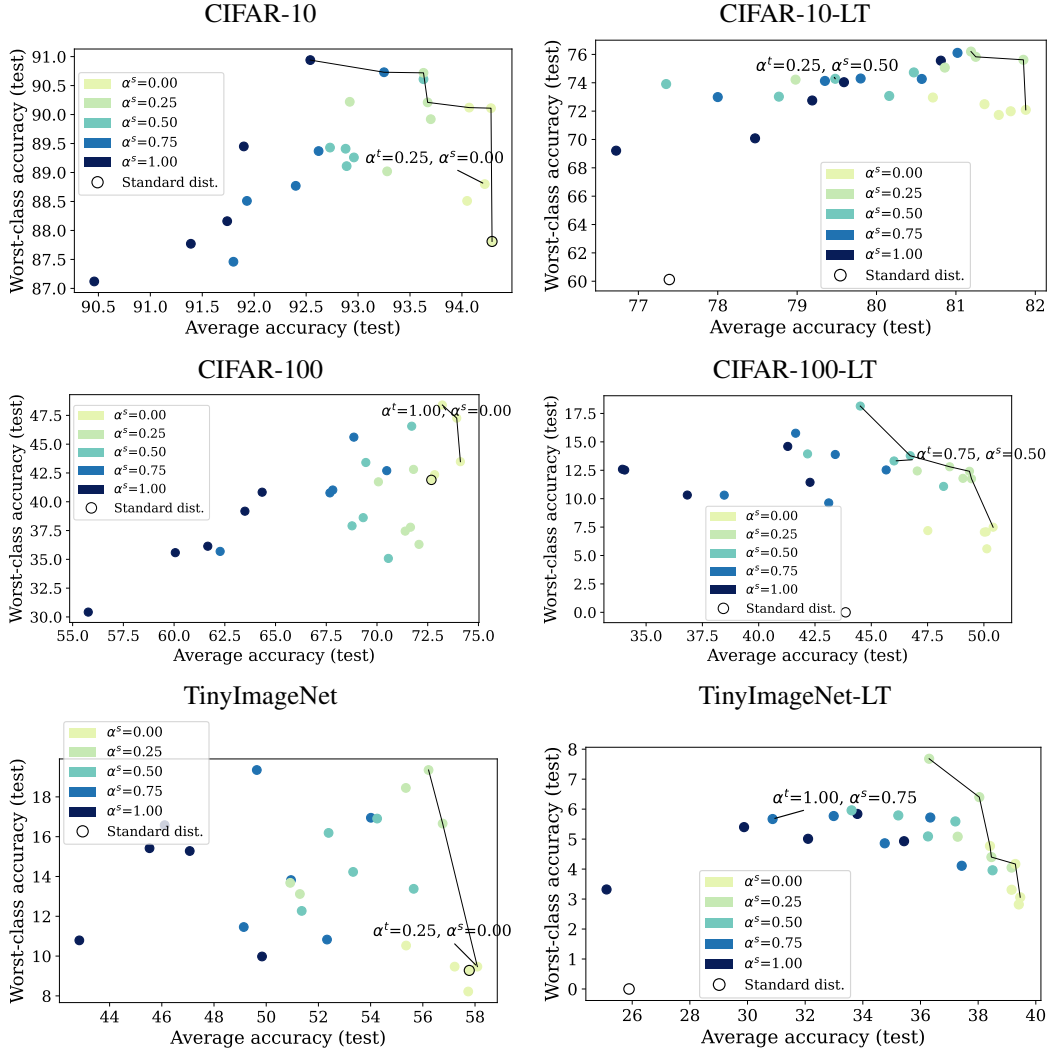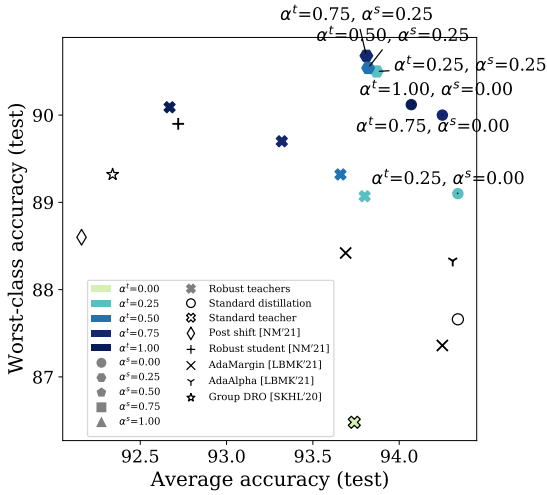
Figure 1: All $\alpha^t, \alpha^s$ combinations for all datasets on test. The black line traces out the Pareto frontier. Average accuracy is roughly determined by $\alpha^s$. The labeled point corresponds to the "best" combination selected in Table 2 based on validation criteria, but other domain-specific trade-off criteria could yield any of these other points.

robust $L^{\text{tdf}}$ combinations also Pareto dominated most of the baselines that all used the standard teacher. Together, these results highlight the fact that in robust distillation, the teacher's training objective is important and should be tailored to the desired final accuracy/robustness trade-off (perhaps using a held-out validation sample with some domain-specific criteria in practice). Figure 4 confirms that these results also hold up in a compression setting, where the compressed models can actually even beat their larger teachers.

## F.4 DIFFERENT TEACHERS ON REPEAT TRAININGS

Distillation experimental results in the main paper use the same teacher for all repeat trainings of the student. This captures the variance in the student training process while omitting the variance in the teacher training process. To capture the variance in the full training pipeline, we ran an additional set of experiments where students were trained on different retrained teachers, rather than on the same teacher. We report results on all CIFAR datasets in Table 3. The best teacher/student combinations are identical for all datasets except for CIFAR-10-LT, for which the best teacher/student combinations from Table 3 and Table 1 were both a robust student trained with a balanced teacher, and only differed in whether the validation set contained teacher labels or one-hot labels ($L^{\text{bal}}/L^{\text{rob-d}}$ (one-hot val) in Table 3 vs. $L^{\text{bal}}/L^{\text{rob-d}}$ (teacher val) in Table 1). Note that the first and second rows of Table 1 are already averaged over $m$ retrained teachers ($m = 5$ for TinyImageNet, or

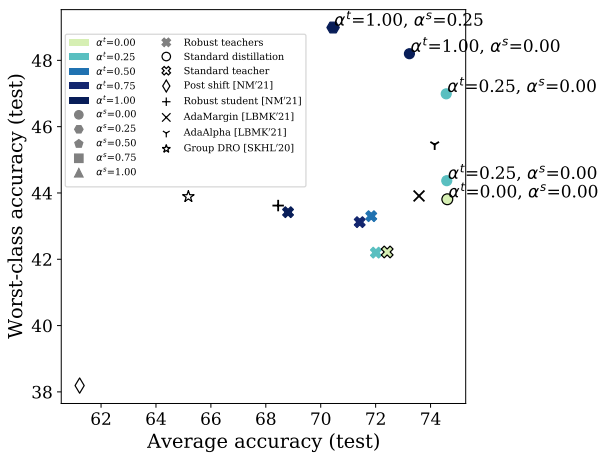**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.75 | 0.25 | $93.81 \pm 0.07$ | $90.68 \pm 0.20$ |
| 0.50 | 0.25 | $93.82 \pm 0.09$ | $90.54 \pm 0.22$ |
| 0.25 | 0.25 | $93.87 \pm 0.08$ | $90.50 \pm 0.18$ |
| 1.00 | 0.00 | $94.07 \pm 0.07$ | $90.12 \pm 0.23$ |
| 0.75 | 0.00 | $94.25 \pm 0.05$ | $90.00 \pm 0.17$ |
| 0.25 | 0.00 | $94.34 \pm 0.06$ | $89.10 \pm 0.31$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $94.34 \pm 0.07$ | $87.66 \pm 0.40$ |
| Post shift [NM'21] | $92.16 \pm 0.18$ | $88.60 \pm 0.35$ |
| Robust student [NM'21] | $92.72 \pm 0.05$ | $89.90 \pm 0.21$ |
| AdaMargin [LBMK'21] | $93.69 \pm 0.06$ | $88.42 \pm 0.36$ |
| AdaAlpha [LBMK'21] | $94.31 \pm 0.01$ | $88.33 \pm 0.14$ |
| Group DRO [SKHL'20] | $92.34 \pm 0.07$ | $89.32 \pm 0.21$ |

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 1.00 | 0.25 | $70.45 \pm 0.16$ | $48.99 \pm 0.72$ |
| 1.00 | 0.00 | $73.23 \pm 0.07$ | $48.20 \pm 1.15$ |
| 0.25 | 0.00 | $74.57 \pm 0.12$ | $46.99 \pm 1.09$ |
| 0.25 | 0.00 | $74.59 \pm 0.09$ | $44.37 \pm 0.58$ |
| 0.00 | 0.00 | $74.61 \pm 0.15$ | $43.81 \pm 0.58$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $74.61 \pm 0.15$ | $43.81 \pm 0.58$ |
| Post shift [NM'21] | $61.22 \pm 0.36$ | $38.19 \pm 0.40$ |
| Robust student [NM'21] | $68.45 \pm 0.13$ | $43.62 \pm 1.27$ |
| AdaMargin [LBMK'21] | $73.58 \pm 0.11$ | $43.91 \pm 1.11$ |
| AdaAlpha [LBMK'21] | $74.15 \pm 0.08$ | $45.46 \pm 0.67$ |
| Group DRO [SKHL'20] | $65.18 \pm 0.08$ | $43.89 \pm 1.12$ |

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.50 | 0.75 | $51.88 \pm 0.18$ | $19.29 \pm 1.27$ |
| 0.75 | 0.50 | $53.60 \pm 0.31$ | $18.98 \pm 0.86$ |
| 0.25 | 0.25 | $56.99 \pm 0.14$ | $18.83 \pm 0.85$ |
| 0.00 | 0.25 | $57.26 \pm 0.15$ | $14.44 \pm 0.91$ |
| 0.75 | 0.00 | $57.35 \pm 0.17$ | $9.47 \pm 1.76$ |
| 0.50 | 0.00 | $57.74 \pm 0.20$ | $8.22 \pm 1.09$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $57.83 \pm 0.13$ | $6.32 \pm 2.31$ |
| Post shift [NM'21] | $43.02 \pm 0.79$ | $14.39 \pm 1.13$ |
| Robust student [NM'21] | $48.06 \pm 0.24$ | $16.27 \pm 0.43$ |
| AdaMargin [LBMK'21] | $52.45 \pm 0.08$ | $15.41 \pm 0.71$ |
| AdaAlpha [LBMK'21] | $57.22 \pm 0.08$ | $7.62 \pm 2.17$ |
| Group DRO [SKHL'20] | $48.78 \pm 0.21$ | $11.38 \pm 1.79$ |

Figure 2: Trade-offs in worst-class test accuracy vs. average test accuracy for CIFAR-10 and CIFAR-100 distilling from ResNet-56 to ResNet-56, and TinyImageNet distilling from ResNet-18 to ResNet-18. All baseline results that require a teacher use the "standard teacher" (trained using $L^{\text{std}}$), as done in the original papers. For methods run multiple times with multiple hyperparameters (e.g. temperatures), all Pareto efficient results are shown in the plot, but the tables show only the baseline results with the best worst-class accuracy (on the validation set). The highlighted row indicates the model with the highest worst-class accuracy that also achieves at least as high average accuracy as *standard distillation*.
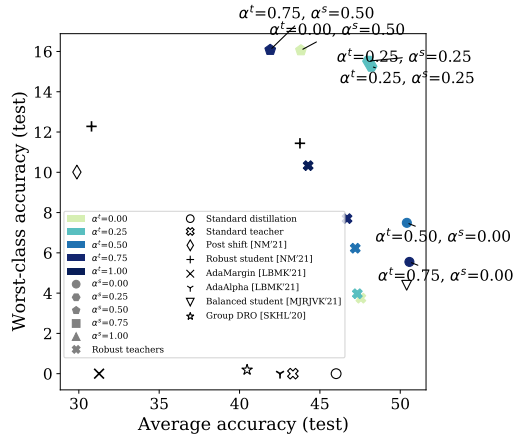
**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.75 | 0.75 | $80.86 \pm 0.09$ | $75.58 \pm 0.17$ |
| 0.75 | 0.50 | $81.12 \pm 0.11$ | $75.52 \pm 0.22$ |
| 0.00 | 0.75 | $81.40 \pm 0.10$ | $75.15 \pm 0.38$ |
| 0.00 | 0.50 | $81.82 \pm 0.11$ | $75.13 \pm 0.24$ |
| 0.00 | 0.25 | $81.89 \pm 0.08$ | $73.09 \pm 0.32$ |
| 0.00 | 0.00 | $81.94 \pm 0.16$ | $70.61 \pm 0.39$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $77.39 \pm 0.10$ | $60.12 \pm 0.56$ |
| Post shift [NM'21] | $78.28 \pm 0.05$ | $74.33 \pm 0.09$ |
| Robust student [NM'21] | $80.05 \pm 0.13$ | $74.91 \pm 0.24$ |
| AdaMargin [LBMK'21] | $72.69 \pm 0.24$ | $47.52 \pm 0.95$ |
| AdaAlpha [LBMK'21] | $70.83 \pm 0.28$ | $43.64 \pm 1.09$ |
| Group DRO [SKHL'20] | $74.39 \pm 0.17$ | $59.93 \pm 0.59$ |

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.75 | 0.50 | $41.91 \pm 0.15$ | $16.08 \pm 0.52$ |
| 0.00 | 0.50 | $43.82 \pm 0.14$ | $16.06 \pm 0.89$ |
| 0.25 | 0.25 | $48.01 \pm 0.09$ | $15.52 \pm 0.41$ |
| 0.25 | 0.25 | $48.20 \pm 0.11$ | $15.26 \pm 0.73$ |
| 0.50 | 0.00 | $50.41 \pm 0.11$ | $7.49 \pm 0.72$ |
| 0.75 | 0.00 | $50.57 \pm 0.18$ | $5.55 \pm 0.54$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $46.01 \pm 0.16$ | $0.00 \pm 0.00$ |
| Post shift [NM'21] | $29.88 \pm 0.61$ | $10.01 \pm 0.72$ |
| Robust student [NM'21] | $30.79 \pm 0.18$ | $12.28 \pm 0.46$ |
| AdaMargin [LBMK'21] | $31.26 \pm 0.21$ | $0.00 \pm 0.00$ |
| AdaAlpha [LBMK'21] | $42.52 \pm 0.08$ | $0.00 \pm 0.00$ |
| Balanced student [MJRJVK'21] | $50.40 \pm 0.12$ | $4.39 \pm 0.66$ |
| Group DRO [SKHL'20] | $40.47 \pm 0.17$ | $0.19 \pm 0.17$ |

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-10 acc. |
|---|---|---|---|
| 1.00 | 0.25 | $36.28 \pm 0.17$ | $7.98 \pm 0.21$ |
| 0.75 | 0.25 | $37.62 \pm 0.15$ | $6.25 \pm 0.12$ |
| 0.00 | 0.25 | $38.44 \pm 0.13$ | $5.90 \pm 0.45$ |
| 0.50 | 0.00 | $39.29 \pm 0.09$ | $4.17 \pm 0.34$ |
| 0.25 | 0.00 | $39.57 \pm 0.06$ | $3.68 \pm 0.30$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-10 acc. |
|---|---|---|
| Standard distillation | $26.05 \pm 0.18$ | $0.00 \pm 0.00$ |
| Post shift [NM'21] | $21.32 \pm 0.49$ | $2.58 \pm 0.42$ |
| Robust student [NM'21] | $21.59 \pm 0.19$ | $1.55 \pm 0.37$ |
| AdaMargin [LBMK'21] | $4.41 \pm 0.09$ | $0.00 \pm 0.00$ |
| AdaAlpha [LBMK'21] | $27.95 \pm 0.14$ | $0.00 \pm 0.00$ |
| Balanced student [MJRJVK'21] | $30.43 \pm 0.06$ | $0.20 \pm 0.18$ |
| Group DRO [SKHL'20] | $27.78 \pm 0.13$ | $0.00 \pm 0.00$ |

Figure 3: Trade-offs in worst-class test accuracy vs. average test accuracy for CIFAR-10-LT, CIFAR-100-LT, and TinyImageNet-LT under self-distillation. All baseline results that require a teacher use the "standard teacher" (trained using $L^{\text{std}}$), as done in the original papers. For methods run multiple times with multiple hyperparameters (e.g. temperatures), all Pareto efficient results are shown in the plot, but the tables show only the baseline results with the best worst-class accuracy (on the validation set). The highlighted row indicates the model with the highest worst-class (or worst-10) accuracy that also achieves at least as high average accuracy as *standard distillation* (within error margins). Note that the for the LT datasets, $L^{\text{tdf}}$ mixes between $L^{\text{bal}}$ and $L^{\text{rob}}$.
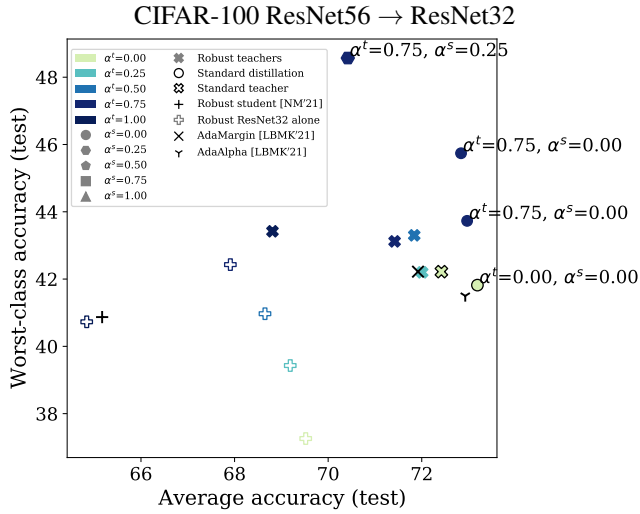
**CIFAR-10 ResNet56 → ResNet32**

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.00 | 0.25 | $93.08 \pm 0.07$ | $89.85 \pm 0.22$ |
| 1.00 | 0.00 | $93.38 \pm 0.05$ | $89.56 \pm 0.20$ |
| 0.75 | 0.00 | $93.58 \pm 0.09$ | $88.91 \pm 0.25$ |
| 1.00 | 0.00 | $93.59 \pm 0.06$ | $88.88 \pm 0.36$ |
| 0.75 | 0.00 | $93.61 \pm 0.05$ | $88.44 \pm 0.33$ |
| 0.25 | 0.00 | $93.74 \pm 0.07$ | $88.41 \pm 0.32$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $93.71 \pm 0.05$ | $86.98 \pm 0.36$ |
| Robust student [NM'21] | $91.57 \pm 0.08$ | $88.57 \pm 0.18$ |
| AdaMargin [LBMK'21] | $92.09 \pm 0.09$ | $83.57 \pm 0.64$ |
| AdaAlpha [LBMK'21] | $93.52 \pm 0.11$ | $85.41 \pm 0.45$ |



**CIFAR-100 ResNet56 → ResNet32**

**Pareto efficient robust distillation results (test)**

| $\alpha^t$ | $\alpha^s$ | Average acc. | Worst-class acc. |
|---|---|---|---|
| 0.75 | 0.25 | $70.42 \pm 0.14$ | $48.57 \pm 0.55$ |
| 0.75 | 0.00 | $72.84 \pm 0.22$ | $45.74 \pm 1.57$ |
| 0.75 | 0.00 | $72.97 \pm 0.18$ | $43.73 \pm 1.72$ |

**Baseline results (test)**

| Baseline | Average acc. | Worst-class acc. |
|---|---|---|
| Standard distillation | $73.19 \pm 0.10$ | $41.82 \pm 1.12$ |
| Robust student [NM'21] | $65.17 \pm 0.11$ | $40.87 \pm 0.89$ |
| AdaMargin [LBMK'21] | $71.92 \pm 0.17$ | $42.22 \pm 1.65$ |
| AdaAlpha [LBMK'21] | $72.93 \pm 0.09$ | $41.50 \pm 1.14$ |

Figure 4: Trade-offs in worst-class test accuracy vs. average test accuracy for CIFAR-10 and CIFAR-100 distilling from ResNet-56 to ResNet-32. All baseline results that require a teacher use the "standard teacher" (trained using $L^{\text{std}}$), as done in the original papers. For methods run multiple times with multiple hyperparameters (e.g. temperatures), all Pareto efficient results are shown in the plot, but the tables show only the baseline results with the best worst-class accuracy (on the validation set). The highlighted row indicates the model with the highest worst-class accuracy that also achieves at least as high average accuracy as *standard distillation* (within error margins).

$m = 10$ for CIFAR datasets), and those same $m$ teachers are used in the repeat trainings in Table 3.

Table 3: Comparison using different teachers for student retrainings for self-distilled teacher/student combos on test. For each student/teacher objective pair, we train $m = 10$ students total on each of $m = 10$ distinct retrained teachers. For comparability, the same set of $m$ teachers is used for each student. This differs from Table 1 in that in Table 1, the students are retrained on each repeat using the same teacher (arbitrarily selected). Otherwise, setups are the same as in Table 1.

| | | **CIFAR-10** Teacher Obj. | | **CIFAR-100** Teacher Obj. | |
|---|---|---|---|---|---|
| | | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ |
| | $L^{\text{std-d}}$ | $87.09 \pm 0.51$ | $89.68 \pm 0.20$ | $44.21 \pm 0.57$ | $\mathbf{47.79 \pm 0.82}$ |
| | | $(93.78 \pm 0.22)$ | $(93.74 \pm 0.07)$ | $74.6 \pm 0.11$ | $73.48 \pm 0.11$ |
| Student Obj. | $L^{\text{rob-d}}$ (teacher val) | $\mathbf{90.62 \pm 0.19}$ | $87.12 \pm 0.38$ | $39.7 \pm 1.32$ | $31.09 \pm 1.21$ |
| | | $(92.58 \pm 0.08)$ | $(90.46 \pm 0.08)$ | $(64.28 \pm 0.41)$ | $(55.39 \pm 0.28)$ |
| | $L^{\text{rob-d}}$ (one-hot val) | $88.15 \pm 0.66$ | $86.44 \pm 0.52$ | $39.44 \pm 0.94$ | $39.65 \pm 0.59$ |
| | | $(91.03 \pm 0.47)$ | $(90.16 \pm 0.42)$ | $(61.23 \pm 0.36)$ | $(60.89 \pm 0.29)$ |

| | | **CIFAR-10-LT** Teacher Obj. | | | **CIFAR-100-LT** Teacher Obj. | | |
|---|---|---|---|---|---|---|---|
| | | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ |
| | $L^{\text{std-d}}$ | $60.12 \pm 0.56$ | $66.13 \pm 0.47$ | $69.75 \pm 0.52$ | $0.00 \pm 0.00$ | $1.41 \pm 0.41$ | $9.17 \pm 0.74$ |
| | | $(77.39 \pm 0.10)$ | $(79.16 \pm 0.20)$ | $(80.73 \pm 0.08)$ | $(45.84 \pm 0.13)$ | $(49.67 \pm 0.20)$ | $(48.55 \pm 0.14)$ |
| Student Obj. | $L^{\text{bal-d}}$ | $72.41 \pm 0.52$ | $71.49 \pm 0.30$ | $71.70 \pm 0.33$ | $5.83 \pm 0.54$ | $5.94 \pm 0.50$ | $8.37 \pm 0.72$ |
| | | $(81.97 \pm 0.11)$ | $(81.20 \pm 0.15)$ | $(80.29 \pm 0.11)$ | $(50.58 \pm 0.15)$ | $(50.85 \pm 0.14)$ | $(48.16 \pm 0.20)$ |
| | $L^{\text{rob-d}}$ (teacher val) | $62.77 \pm 0.58$ | $73.09 \pm 0.34$ | $68.04 \pm 0.47$ | $10.53 \pm 0.76$ | $12.04 \pm 0.89$ | $9.66 \pm 1.15$ |
| | | $(77.18 \pm 0.15)$ | $(80.03 \pm 0.22)$ | $(75.36 \pm 0.25)$ | $(33.69 \pm 0.14)$ | $(34.08 \pm 0.12)$ | $(37.10 \pm 0.15)$ |
| | $L^{\text{rob-d}}$ (one-hot val) | $\mathbf{75.10 \pm 0.36}$ | $\mathbf{75.10 \pm 0.50}$ | $74.16 \pm 0.34$ | $10.74 \pm 0.44$ | $11.95 \pm 0.69$ | $\mathbf{12.87 \pm 0.81}$ |
| | | $(79.27 \pm 0.13)$ | $(79.07 \pm 0.20)$ | $(78.11 \pm 0.14)$ | $(30.36 \pm 0.39)$ | $(31.00 \pm 0.16)$ | $(31.62 \pm 0.34)$ |

## F.5 ADAALPHA AND ADAMARGIN COMPARISONS WITH DIFFERENT TEACHERS

We include and discuss additional comparisons to the AdaMargin and AdaAlpha methods Lukasik et al. [2022], which each define additional ways to modify the student training algorithm (see Section 4). In Table 2, we show results with each of these methods using the standard teacher, as done in the original paper. However, in this section we extend these results by also applying AdaMargin and AdaAlpha with different teachers trained with the robust and balanced objectives. Table 4 compares the results of AdaMargin and AdaAlpha for these different teachers under the same self distillation setup as Table 1.

Overall, the use of a robust teacher leads to marked improvements for students trained by AdaMargin and AdaAlpha. For the balanced datasets, AdaMargin was competitive with the robust and standard students: on CIFAR-100 and TinyImageNet, AdaMargin combined with the robust teacher and the standard teacher (respectively) achieved worst-class accuracies that are statistically comparable to the best worst-class accuracies in Table 1. However, on the long-tailed datasets, AdaAlpha and AdaMargin did not achieve worst-class accuracies as high as other teacher/student combinations. This suggests that the AdaMargin method can work well on balanced datasets in combination with a robust teacher, but other combinations of standard/balanced/robust objectives are valuable for long-tailed datasets.

Relative to each other, AdaMargin usually achieved higher worst-class accuracy than AdaAlpha, whereas AdaAlpha often achieved higher average accuracy.

## F.6 GROUP DRO COMPARISON

Sagawa et al. [2020] propose a group DRO algorithm to improve long tail performance without distillation. In this section we present additional experimental comparisons to Algorithm 1 from Sagawa et al. [2020]. This differs from our robust optimization methodology in Section 3.1 in two key ways: *(i)* we apply a margin-based surrogates of Menon et al. [2021b], and *(ii)* we use a validation set to update the Lagrange multipliers $\lambda$ in Algorithm 2. Table 5 shows results from running group DRO directly as specified in Algorithm 1 in Sagawa et al. [2020], as well as a variant where we use the validation set

Table 4: Results for AdaAlpha and AdaMargin baselines for different teachers under self-distillation. For all CIFAR datasets, self-distillation is done from ResNet56 → ResNet56. For TinyImageNet, self-distillation is done from ResNet18 → ResNet18. Worst-class accuracy shown above (or worst-10 accuracy for TinyImageNet-LT), and average accuracy is shown in parentheses below. The temperature hyperparameter was tuned to maximize worst-class accuracy on the held-out validation set. Mean and standard error are reported over 5 repeats for all datasets.

| | **CIFAR-10** Teacher Obj. | | **CIFAR-100** Teacher Obj. | | **TinyImageNet** Teacher Obj. | |
| | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ |
|---|---|---|---|---|---|---|
| Ada Alpha | $88.33 \pm 0.14$ ($94.31 \pm 0.01$) | $89.96 \pm 0.44$ ($93.97 \pm 0.07$) | $43.50 \pm 0.62$ ($73.96 \pm 0.09$) | $45.59 \pm 0.82$ ($71.42 \pm 0.14$) | $11.11 \pm 1.29$ ($61.13 \pm 0.09$) | $16.58 \pm 1.67$ ($56.84 \pm 0.15$) |
| Ada Margin | $87.36 \pm 0.06$ ($94.25 \pm 0.02$) | $90.37 \pm 0.26$ ($94.02 \pm 0.12$) | $43.91 \pm 1.11$ ($73.58 \pm 0.11$) | $47.78 \pm 0.96$ ($70.92 \pm 0.09$) | $18.17 \pm 3.89$ ($61.3 \pm 0.28$) | $17.84 \pm 1.77$ ($55.77 \pm 0.32$) |

| | **CIFAR-10-LT** Teacher Obj. | | | **CIFAR-100-LT** Teacher Obj. | | |
| | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ |
|---|---|---|---|---|---|---|
| Ada Alpha | $41.90 \pm 0.44$ ($71.67 \pm 0.08$) | $66.23 \pm 0.39$ ($77.87 \pm 0.16$) | $71.17 \pm 0.32$ ($79.66 \pm 0.13$) | $0.00 \pm 0.00$ ($42.52 \pm 0.08$) | $1.46 \pm 0.61$ ($45.44 \pm 0.14$) | $9.15 \pm 0.54$ ($45.64 \pm 0.11$) |
| Ada Margin | $47.52 \pm 0.95$ ($72.69 \pm 0.24$) | $66.74 \pm 0.35$ ($78.20 \pm 0.09$) | $70.33 \pm 0.50$ ($78.87 \pm 0.12$) | $0.00 \pm 0.00$ ($31.26 \pm 0.21$) | $0.00 \pm 0.00$ ($34.06 \pm 0.12$) | $12.46 \pm 0.36$ ($42.90 \pm 0.07$) |

| | **TinyImageNet-LT** Teacher Obj. | | |
| | $L^{\text{std}}$ | $L^{\text{bal}}$ | $L^{\text{rob}}$ |
|---|---|---|---|
| Ada Alpha | $0.00 \pm 0.00$ ($28.14 \pm 0.12$) | $0.00 \pm 0.00$ ($0.50 \pm 0.00$) | $0.00 \pm 0.00$ ($0.50 \pm 0.00$) |
| Ada Margin | $0.00 \pm 0.00$ ($9.18 \pm 0.09$) | $0.00 \pm 0.00$ ($7.92 \pm 0.10$) | $0.41 \pm 0.17$ ($23.08 \pm 0.15$) |

to update Lagrange multipliers in group DRO (labeled as "with vali" in Table 5). Table 5 shows that this latter variant "with vali" performs better than the original version without a validation set; thus, for the results in Figures 2 and 3, we report these better results marked in Table 5 as "with vali." Overall, this comparison shows that $L^{\text{rob}}$ is comparable to group DRO, and that robust distillation protocols can outperform group DRO alone.

Table 5: Results from comparison to group DRO (Algorithm 1 in Sagawa et al. [2020]) without distillation. "No vali" uses the training set to update group Lagrange multipliers, as done originally by Sagawa et al. [2020]. "With vali" uses the validation set to compute group Lagrange multipliers as done in all other experiments in our paper. Worst-class accuracy is shown above, and balanced accuracy is shown in parentheses below. Mean and standard error are shown over 5 repeats.

| **CIFAR-10** group DRO | | **CIFAR-100** group DRO | | **TinyImageNet** group DRO | |
| No vali | With vali | No vali | With vali | No vali | With vali |
|---|---|---|---|---|---|
| $86.65 \pm 0.49$ ($93.61 \pm 0.09$) | $89.32 \pm 0.21$ ($92.34 \pm 0.07$) | $40.35 \pm 1.18$ ($70.25 \pm 0.17$) | $43.89 \pm 1.12$ ($65.18 \pm 0.08$) | $0.00 \pm 0.00$ ($6.55 \pm 0.41$) | $9.17 \pm 1.55$ ($47.67 \pm 0.22$) |

| **CIFAR-10-LT** group DRO | | **CIFAR-100-LT** group DRO | | **TinyImageNet-LT** group DRO | |
| No vali | With vali | No vali | With vali | No vali | With vali |
|---|---|---|---|---|---|
| $51.59 \pm 2.49$ ($71.94 \pm 0.75$) | $59.93 \pm 0.59$ ($74.39 \pm 0.17$) | $0.00 \pm 0.00$ ($39.81 \pm 0.23$) | $0.19 \pm 0.17$ ($40.47 \pm 0.17$) | $0.00 \pm 0.00$ ($9.79 \pm 0.40$) | $0.00 \pm 0.00$ ($22.49 \pm 0.10$) |

## F.7 ADDITIONAL IMAGENET COMPARISONS

Here we present additional results when training ResNet-18 teachers and students on ImageNet. Table 6 includes measures of worst-1 accuracy, worst-10 accuracy, worst-50 accuracy, and worst-100 accuracy.

Table 6: ImageNet comparison of ResNet-18 teacher/student combos on test. Average worst-1/10/100 accuracy shown above, standard accuracy shown in parentheses below. The combination with the best worst-class accuracy is bolded. Mean and standard error are reported over up to 5 repeats.

| | | Worst-1 Accuracy | | Worst-10 Accuracy | | Worst-50 Accuracy | | Worst-100 Accuracy | |
|---|---|---|---|---|---|---|---|---|---|
| | | **ImageNet** Teacher Obj. | | **ImageNet** Teacher Obj. | | **ImageNet** Teacher Obj. | | **ImageNet** Teacher Obj. | |
| | | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ | $L^{\text{std}}$ | $L^{\text{rob}}$ |
| Student Obj. | none | 0.00 (67.29) | 10.71 (63.10) | 11.54 (67.29) | 17.13 (63.10) | 24.01 (67.29) | 25.08 (63.10) | 30.35 (67.29) | 29.96 (63.10) |
| | Post shift | 8.70 (48.62) | 3.57 (48.83) | 16.15 (48.62) | 11.64 (48.83) | 21.85 (48.62) | 18.86 (48.83) | 25.58 (48.62) | 23.17 (48.83) |
| | $L^{\text{std-d}}$ | 3.20 ± 1.33 (65.46 ± 0.05) | 3.79 ± 0.11 (64.54 ± 0.01) | 10.07 ± 0.27 (65.46 ± 0.05) | 10.22 ± 0.33 (64.54 ± 0.01) | 20.30 ± 0.39 (65.46 ± 0.05) | 22.61 ± 0.32 (64.54 ± 0.01) | 26.45 ± 0.25 (65.46 ± 0.05) | 29.01 ± 0.21 (64.54 ± 0.01) |
| | $L^{\text{rob-d}}$ (teacher val) | 0.00 ± 0.00 (59.60 ± 0.10) | 0.00 ± 0.00 (51.01 ± 0.12) | 1.18 ± 0.02 (59.60 ± 0.10) | 1.47 ± 0.04 (51.01 ± 0.12) | 13.02 ± 0.16 (59.60 ± 0.10) | 6.85 ± 0.13 (51.01 ± 0.12) | 21.00 ± 0.19 (59.60 ± 0.10) | 11.26 ± 0.16 (51.01 ± 0.12) |
| | $L^{\text{rob-d}}$ (one-hot val) | 0.00 ± 0.00 (59.65 ± 0.01) | 0.00 ± 0.00 (55.34 ± 0.00) | 8.32 ± 1.04 (59.65 ± 0.01) | 5.99 ± 0.00 (55.34 ± 0.00) | 18.77 ± 0.03 (59.65 ± 0.01) | 16.82 ± 0.00 (55.34 ± 0.00) | 23.95 ± 0.31 (59.65 ± 0.01) | 22.16 ± 0.00 (55.34 ± 0.00) |

## References

Andrew Cotter, Heinrich Jiang, Serena Wang, Taman Narayan, Seungil You, Karthik Sridharan, and Maya R. Gupta. Optimization with non-differentiable constraints with applications to fairness, recall, churn, and other goals. *Journal of Machine Learning Research (JMLR)*, 20(172):1–59, 2019.

Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9268–9277, 2019.

Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 212–220, 2017.

Michal Lukasik, Srinadh Bhojanapalli, Aditya Krishna Menon, and Sanjiv Kumar. Teacher's pet: understanding and mitigating biases in distillation. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=ph3AYXpwEb.

Aditya K Menon, Ankit Singh Rawat, Sashank Reddi, Seungyeon Kim, and Sanjiv Kumar. A statistical perspective on distillation. In *International Conference on Machine Learning*, pages 7632–7642. PMLR, 2021a.

Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. In *International Conference on Learning Representations (ICLR)*, 2021b.

Harikrishna Narasimhan and Aditya K Menon. Training over-parameterized models with non-decomposable objectives. *Advances in Neural Information Processing Systems*, 34, 2021.

Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020.

Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and trends in Machine Learning*, 4 (2):107–194, 2011.

Maurice Sion. On general minimax theorems. *Pacific Journal of mathematics*, 8(1):171–176, 1958.

Robert C Williamson, Elodie Vernet, and Mark D Reid. Composite multiclass losses. *Journal of Machine Learning Research*, 17:1–52, 2016.