# Appendix
# Implementation Details

## Hyperparameters

Table 2: Hyperparameter settings, which are inherited from the implementation by Oliver et al. (Oliver et al. 2018), where hyperparameters are tuned on a validation set. An Adam optimiser is deployed with the same learning rate decay schedule. The following hyperparameters are used for all experiments. Note: we adopt the *same* ramp-up function for all methods.

| shared | |
| --- | --- |
| training iterations | 500,000 |
| coefficient rampup from 0 util | 200,000 |
| learning decay factor | 0.2 |
| learning decay at iteration | 400,000 |
| **supervised baseline** | |
| initial learning rate | 0.003 |
| **pseudo-label** | |
| initial learning rate | 0.003 |
| max consistency coefficient | 1 |
| pseudo-label threshold | 0.95 |
| **VAT** | |
| initial learning rate | 0.003 |
| max consistency coefficient | 0.3 |
| VAT $\epsilon, \xi$ | $6.0, 10^{-6}$ |
| **$\Pi$-Model** | |
| initial learning rate | 0.0003 |
| **Mean-Teacher** | |
| initial learning rate | 0.0004 |
| max consistency coefficient | 8 |
| Exponential moving average decay | 0.95 |
| **SWA** | |
| initial learning rate | 0.001 |
| max consistency coefficient | 8 |
| weight averaging interval | 5,000 |
| **UASD** | |
| initial learning rate | 0.001 |
| max distillation coefficient | 1 |

## Data Augmentation

Table 3: Data augmentation. Note: ZCA image pre-processing is *only* applied on CIFAR10.

| dataset | gaussian noise $\sigma = 0.15$ | horizontal flip $p = 0.5$ | random translation $[-2, +2]$ |
| --- | --- | --- | --- |
| CIFAR10 | ✓ | ✓ | ✓ |
| CIFAR100 | ✗ | ✓ | ✓ |
| TinyImageNet | ✗ | ✓ | ✓ |

# Evaluation Protocols

Table 4: Evaluation protocols of SSL under class mismatch. $p\%$: Class distribution mismatch proportion among unlabelled data. $K$: number of known classes in labelled data. $L_{num}$: Labels per class.

| Dataset | $p\%$ | $K$ | $L_{num}$ | Labelled classes | Unlabelled classes |
|---|---|---|---|---|---|
| CIFAR10 | 0<br>25<br>50<br>75<br>100 | 400 | 6 | 2,3,4,5,6,7 | 4,5,6,7<br>0,5,6,7<br>0,1,6,7<br>0,1,8,7<br>0,1,8,9 |
| CIFAR100 | 50 | 50 | 100 | 0-50 | 25-75 |
| TinyImageNet | 50 | 100 | 100 | 0-100 | 50-150 |
| CIFAR100 + TinyImageNet | 86.5 | 100 | 100 | CIFAR100 | TinyImageNet |

# Tabular Results

## Evaluation on CIFAR10

Table 5: Evaluation under varying class distribution mismatch proportion on CIFAR10. Test error rates are reported at the point of lowest validation error. Results with reduction in error rate compared to supervised learning baseline are highlighted in **bold**. Best results are highlighted in <span style="color:red">**red**</span>.

| Method | Class Distribution Mismatch Proportion | | | | |
|---|---|---|---|---|---|
|  | 0% | 25% | 25% | 50% | 100% |
| baseline |  |  | $24.03 \pm 0.75$ |  |  |
| pseudo-label | $\mathbf{22.37 \pm 0.66}$ | $\mathbf{24.02 \pm 0.86}$ | $25.37 \pm 0.86$ | $26.11 \pm 1.91$ | $26.29 \pm 0.71$ |
| VAT | $\mathbf{20.63 \pm 1.77}$ | $\mathbf{23.08 \pm 0.49}$ | $\mathbf{23.78 \pm 0.70}$ | $25.52 \pm 0.84$ | $26.23 \pm 0.37$ |
| Π-Model | $\mathbf{21.56 \pm 1.29}$ | $24.80 \pm 1.32$ | $25.92 \pm 1.61$ | $26.43 \pm 0.81$ | $26.61 \pm 0.79$ |
| Temporal Ensembling | $\mathbf{21.93 \pm 0.43}$ | $24.23 \pm 0.96$ | $25.66 \pm 1.21$ | $26.33 \pm 0.56$ | $27.00 \pm 1.39$ |
| Mean-Teacher | $\mathbf{21.68 \pm 0.88}$ | $24.13 \pm 1.22$ | $24.79 \pm 1.53$ | $25.90 \pm 1.00$ | $26.78 \pm 0.38$ |
| SWA | $\mathbf{21.63 \pm 0.38}$ | $\mathbf{23.31 \pm 0.85}$ | $\mathbf{23.70 \pm 0.61}$ | $\mathbf{23.90 \pm 0.85}$ | $24.11 \pm 0.65$ |
| UASD (ours) | <span style="color:red">$\mathbf{20.59 \pm 0.51}$</span> | <span style="color:red">$\mathbf{21.34 \pm 0.52}$</span> | <span style="color:red">$\mathbf{21.88 \pm 0.69}$</span> | <span style="color:red">$\mathbf{22.39 \pm 0.48}$</span> | <span style="color:red">$\mathbf{22.49 \pm 0.90}$</span> |

Table 6: Evaluation under varying class distribution mismatch proportion on CIFAR10. Test error rates are reported as the median of last 20 epochs. Results with reduction in error rate compared to supervised learning baseline are highlighted in **bold**. Best results are highlighted in <span style="color:red">**red**</span>.

| Method | Class Distribution Mismatch Proportion | | | | |
|---|---|---|---|---|---|
|  | 0% | 25% | 25% | 50% | 100% |
| baseline |  |  | $23.82 \pm 0.61$ |  |  |
| pseudo-label | $\mathbf{22.70 \pm 0.42}$ | $24.42 \pm 0.87$ | $26.47 \pm 1.01$ | $27.82 \pm 1.10$ | $28.07 \pm 1.01$ |
| VAT | $\mathbf{23.07 \pm 0.49}$ | $27.27 \pm 1.36$ | $27.45 \pm 2.17$ | $28.46 \pm 2.62$ | $28.79 \pm 1.11$ |
| Π-Model | $\mathbf{22.97 \pm 0.46}$ | $26.48 \pm 0.66$ | $29.01 \pm 2.67$ | $28.19 \pm 0.97$ | $29.43 \pm 1.88$ |
| Temporal Ensembling | $\mathbf{22.45 \pm 0.59}$ | $25.33 \pm 0.81$ | $26.94 \pm 0.57$ | $27.59 \pm 0.62$ | $28.16 \pm 0.70$ |
| Mean-Teacher | $\mathbf{22.09 \pm 0.57}$ | $25.40 \pm 0.41$ | $26.46 \pm 0.78$ | $27.83 \pm 1.43$ | $29.09 \pm 1.44$ |
| SWA | $\mathbf{21.70 \pm 0.34}$ | $\mathbf{23.36 \pm 0.74}$ | $23.83 \pm 0.61$ | $24.15 \pm 0.90$ | $24.31 \pm 0.55$ |
| UASD (ours) | <span style="color:red">$\mathbf{20.55 \pm 0.41}$</span> | <span style="color:red">$\mathbf{21.66 \pm 0.71}$</span> | <span style="color:red">$\mathbf{22.01 \pm 0.78}$</span> | <span style="color:red">$\mathbf{22.31 \pm 0.65}$</span> | <span style="color:red">$\mathbf{22.63 \pm 0.78}$</span> |

## Ablative Analysis

Table 7: Evaluation under varying class distribution mismatch proportion on CIFAR10. "size": ensemble size. Test error rates are reported as the median of last 20 epochs. Results with reduction in error rate compared to supervised learning baseline are highlighted in **bold**. Best results are highlighted in **red**.

| Method | Class Distribution Mismatch Proportion | | | | |
|---|---|---|---|---|---|
|  | 0% | 25% | 25% | 50% | 100% |
| baseline |  |  | $23.82 \pm 0.61$ |  |  |
| size 10+ | $\mathbf{21.12 \pm 0.69}$ | $\mathbf{22.64 \pm 0.77}$ | $\mathbf{23.41 \pm 0.56}$ | $24.09 \pm 0.72$ | $24.68 \pm 0.70$ |
| size 100+ | $\mathbf{20.57 \pm 0.39}$ | $\mathbf{\color{red}21.40 \pm 0.41}$ | $\mathbf{22.59 \pm 0.49}$ | $\mathbf{22.44 \pm 0.70}$ | $\mathbf{22.80 \pm 0.55}$ |
| size 1000+ (ours) | $\mathbf{\color{red}20.55 \pm 0.41}$ | $\mathbf{21.66 \pm 0.71}$ | $\mathbf{\color{red}22.01 \pm 0.78}$ | $\mathbf{\color{red}22.31 \pm 0.65}$ | $\mathbf{\color{red}22.63 \pm 0.78}$ |

Table 8: Evaluation under varying class distribution mismatch proportion on CIFAR10. "w/o both": w/o soft distillation and w/o OOD filter. "w/o soft": w/o soft distillation. "w/o OOD": w/o OOD filter. Test error rates are reported as the median of last 20 epochs. Results with reduction in error rate compared to supervised learning baseline are highlighted in **bold**. Best results are highlighted in **red**.

| Method | Class Distribution Mismatch Proportion | | | | |
|---|---|---|---|---|---|
|  | 0% | 25% | 25% | 50% | 100% |
| baseline |  |  | $23.82 \pm 0.61$ |  |  |
| w/o both | $\mathbf{23.50 \pm 0.86}$ | $24.78 \pm 0.64$ | $25.78 \pm 0.57$ | $26.05 \pm 0.79$ | $27.43 \pm 0.74$ |
| w/o soft | $\mathbf{21.84 \pm 0.53}$ | $\mathbf{23.27 \pm 0.60}$ | $24.67 \pm 0.60$ | $24.67 \pm 0.82$ | $25.52 \pm 0.92$ |
| w/o OOD | $\mathbf{21.19 \pm 0.31}$ | $\mathbf{22.34 \pm 0.52}$ | $\mathbf{22.61 \pm 0.99}$ | $\mathbf{22.68 \pm 0.56}$ | $\mathbf{23.11 \pm 0.70}$ |
| Full UASD (ours) | $\mathbf{\color{red}20.55 \pm 0.41}$ | $\mathbf{\color{red}21.66 \pm 0.71}$ | $\mathbf{\color{red}22.01 \pm 0.78}$ | $\mathbf{\color{red}22.31 \pm 0.65}$ | $\mathbf{\color{red}22.63 \pm 0.78}$ |

## Classes in CIFAR100 and TinyImageNet

Table 9: Classes in CIFAR100. Class labels which overlap with (i.e. same as or similar to) TinyImageNet are highlighted in **bold**. Number of same or similar labels: 19.

| Superclass | Class labels |
|---|---|
| aquatic mammals | beaver, dolphin, otter, seal, whale |
| fish | aquarium fish, flatfish, ray, shark, trout |
| flowers | orchids, poppies, roses, sunflowers, tulips |
| food containers | **bottles**, bowls, cans, cups, **plates** |
| fruit and vegetables | apples, **mushrooms**, **oranges**, pears, sweet peppers |
| household electrical devices | clock, **computer keyboard**, **lamp**, telephone, television |
| household furniture | bed, **chair**, couch, **table**, wardrobe |
| insects | **bee**, beetle, **butterfly**, caterpillar, **cockroach** |
| large carnivores | **bear**, leopard, **lion**, tiger, wolf |
| large man-made outdoor things | **bridge**, castle, house, road, skyscraper |
| large natural outdoor scenes | cloud, forest, mountain, plain, sea |
| large omnivores and herbivores | **camel**, cattle, **chimpanzee**, **elephant**, kangaroo |
| medium-sized mammals | fox, porcupine, possum, raccoon, skunk |
| non-insect invertebrates | crab, **lobster**, **snail**, spider, worm |
| people | baby, boy, girl, man, woman |
| reptiles | crocodile, dinosaur, lizard, snake, turtle |
| small mammals | hamster, mouse, rabbit, shrew, squirrel |
| trees | maple, oak, palm, pine, willow |
| vehicles 1 | bicycle, **bus**, motorcycle, pickup truck, **train** |
| vehicles 2 | lawn-mower, rocket, streetcar, tank, **tractor** |

Table 10: Classes in TinyImageNet. Class labels which overlap with (i.e. same as or similar to) CIFAR100 are highlighted in **bold**. Number of same or similar labels: 27. Class distribution mismatch proportion compared to CIFAR100: $(200 - 17)/200 = \mathbf{86.5}\%$.

| ID | Class labels |
|---|---|
| 1-5 | Egyptian cat, reel, volleyball, **rocking chair**, lemon |
| 6-10 | bullfrog, basketball, cliff, espresso, plunger |
| 11-15 | parking meter, German shepherd, **dining table**, monarch, **brown bear** |
| 16-20 | school bus, pizza, guinea pig, umbrella, organ |
| 21-25 | oboe, maypole, goldfish, potpie, hourglass |
| 26-30 | seashore, **computer keyboard**, **Arabian camel**, ice cream, nail |
| 31-35 | space heater, cardigan, baboon, **snail**, coral reef |
| 25-30 | albatross, **spider web**, sea cucumber, backpack, Labrador retriever |
| 36-40 | pretzel, king penguin, **sulphur butterfly**, tarantula, lesser panda |
| 46-50 | **pop bottle**, banana, sock, **cockroach**, projectile |
| 51-55 | **beer bottle**, mantis, freight car, guacamole, remote control |
| 56-60 | European fire salamander, lakeside, **chimpanzee**, pay-phone, fur coat |
| 61-65 | alp, **lampshade**, torch, abacus, moving van |
| 66-70 | barrel, tabby, goose, koala, **bullet train** |
| 71-75 | CD player, teapot, birdhouse, gazelle, academic gown |
| 76-80 | **tractor**, ladybug, miniskirt, golden retriever, triumphal arch |
| 81-85 | cannon, neck brace, sombrero, gasmask, candle |
| 86-90 | desk, frying pan, **bee**, dam, **spiny lobster** |
| 91-95 | police van, iPod, punching bag, beacon, jellyfish |
| 96-100 | wok, potter's wheel, sandal, **pill bottle**, butcher shop |
| 101-105 | slug, hog, cougar, crane, vestment |
| 106-110 | dragonfly, cash machine, **mushroom**, jinrikisha, water tower |
| 111-115 | chest, snorkel, sunglasses, fly, limousine |
| 116-120 | black stork, dugong, sports car, water jug, **suspension bridge** |
| 121-125 | ox, ice lolly, turnstile, Christmas stocking, broom |
| 126-130 | scorpion, wooden spoon, picket fence, rugby ball, sewing machine |
| 131-135 | **steel arch bridge**, Persian cat, refrigerator, barn, apron |
| 136-140 | Yorkshire terrier, swimming trunks, stopwatch, lawn mower, thatch |
| 141-145 | fountain, black widow, bikini, **plate**, teddy |
| 146-150 | barbershop, confectionery, beach wagon, scoreboard, **orange** |
| 151-155 | flagpole, **American lobster**, **trolleybus**, drumstick, dumbbell |
| 156-160 | brass, bow tie, convertible, bighorn, orangutan |
| 161-165 | American alligator, centipede, syringe, go-kart, brain coral |
| 166-170 | sea slug, cliff dwelling, mashed potato, viaduct, military uniform |
| 171-175 | pomegranate, chain, kimono, comic book, trilobite |
| 176-180 | bison, pole, boa constrictor, poncho, bathtub |
| 181-185 | grasshopper, walking stick, Chihuahua, tailed frog, **lion** |
| 186-190 | altar, obelisk, beaker, bell pepper, bannister |
| 191-195 | bucket, magnetic compass, meat loaf, gondola, standard poodle |
| 196-200 | acorn, lifeboat, binoculars, cauliflower, **African elephant** |

# Confidence Score Estimated by Different SSL Methods

Figure 8: Average confidence score (i.e. maximum class probability) on unlabelled data, estimated by different teaching signals *during training* under varying class distribution mismatch proportion (i.e. 0, 25, 50, 75, 100%) on CIFAR10. Most SSL methods are prone to produce ***overconfident*** teaching signals, regardless of the underlying unlabelled class distribution. This hinders the possibility to be aware of uncertainty, and blindly reinforces the overconfident wrong class assignments on those irrelevant unlabelled samples. In contrast, UASD produces soft teaching signals that encode higher uncertainty, and exhibits different levels of confidence score that are clearly stratified to reflect the underlying class distribution mismatch proportions.