

Adding Seemingly Uninformative Labels Helps in Low Data Regimes

Christos Matsoukas^{1 2 3} Albert Bou I Hernandez¹ Yue Liu^{1 2} Karin Dembrower^{4 5} Gisele Miranda^{1 2}
Emir Konuk^{1 2} Johan Fredin Haslum^{1 2 3} Athanasios Zouzos⁶ Peter Lindholm⁴ Fredrik Strand^{4 6}
Kevin Smith^{1 2}

Abstract

Evidence suggests that networks trained on large datasets generalize well not solely because of the numerous training examples, but also class diversity which encourages learning of enriched features. This raises the question of whether this remains true when data is scarce – *is there an advantage to learning with additional labels in low-data regimes?* In this work, we consider a task that requires difficult-to-obtain expert annotations: tumor segmentation in mammography images. We show that, in low-data settings, performance can be improved by complementing the expert annotations with seemingly uninformative labels from non-expert annotators, turning the task into a multi-class problem. We reveal that these gains increase when less expert data is available, and uncover several interesting properties through further studies. We demonstrate our findings on CSAW-S, a new dataset that we introduce here, and confirm them on two public datasets.

1. Introduction

When abundant training data is available, deep learning approaches have achieved remarkable feats, especially in areas like computer vision and natural language processing. However, even within these well-studied domains, there exist many important problems for which data is scarce.

In medicine, large well-labeled datasets are difficult to come by for a variety of reasons. Often, legal regulations and privacy concerns prohibit the distribution of data. Or there

¹School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden ²Science for Life Laboratory, Stockholm, Sweden ³AstraZeneca, Gothenburg, Sweden ⁴Karolinska Institutet, Stockholm, Sweden ⁵Capio Sankt Göran Hospital, Stockholm, Sweden ⁶Karolinska University Hospital, Stockholm, Sweden. Correspondence to: Christos Matsoukas <matsou@kth.se>.

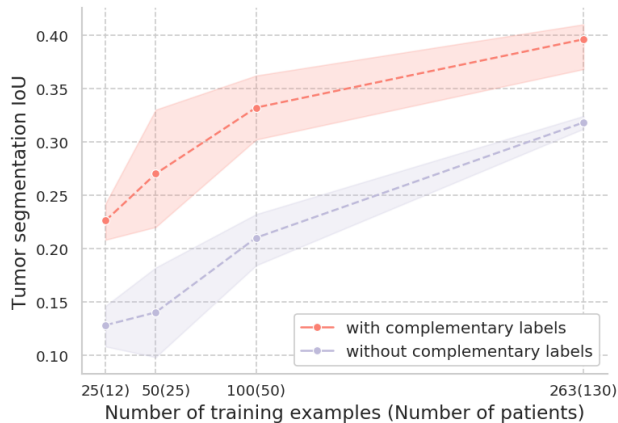


Figure 1. We consider low-data regime tasks where expert data is difficult to obtain, such as tumor segmentation in mammography images from the CSAW-S dataset. A network trained with only expert annotations of the tumors (blue) is significantly outperformed by the same network which is given additional complementary non-expert annotations (red) of breast anatomy (*e.g.* skin, pectoral muscle, nipple), turning the task into a multi-class problem. Performance is measured for varying numbers of training images/patients by intersection over union (IoU) over five repetitions (shaded area indicates 95% CI).

simply may not be enough patients to collect the data from. This may be due to the rare nature of a disease that only affects a limited population, or due to the logistics of collecting data from different administrative regions. Sometimes the techniques used for diagnosis are expensive or require painful procedures (*e.g.* biopsies), and thus data is only collected when deemed absolutely necessary. Finally, there are cases where sufficient data may be available, but crucial expert-level annotations are prohibitively expensive or difficult to obtain. This problem is common in other domains such as astronomy and biology where either data or expert knowledge is difficult to obtain.

In this work, we observe that performance on expert image segmentation tasks can be significantly improved by adding seemingly uninformative annotations during training. Although we cannot offer a definitive explanation for this effect, we surmise that the class diversity encourages learning of enriched features that capture complementary information to the main task. The additional labels do not

contain any direct evidence for the expert task, but they may give indirect support by providing revealing information about other content in the image.

The principal benefit of this approach is that model performance can be increased without changing the architecture or collecting additional expert training examples. The additional labels can be obtained inexpensively and do not need to be of high quality. This allows one to achieve better performance when it is simply not possible to obtain more data, and to optimize the costs of the data acquisition process by considering the trade-offs between collecting new examples, expert annotations, and non-expert annotations.

Many researchers share an intuition that adding information to the learning targets improves performance, and in fact many related techniques discussed in this work are based on this principle (*e.g.* distillation, multi-task learning). Our contribution is to demonstrate and gain insights into this principal for complementary labels, which is not yet well understood.

As our main case study, we introduce CSAW-S, a dataset of mammography images which includes expert annotations of tumors and non-expert annotations of breast anatomy and artifacts in the image (described in Section 4). Using this dataset we experimentally show that tumor segmentation performance, a task that requires expert annotations, is significantly improved by providing additional complementary non-expert labels (*e.g.* skin, nipple, pectoral muscle). We further show that this benefit becomes more prominent as data becomes more scarce.

We validate our findings by demonstrating that the observed effect holds in other domains, using public datasets including CITYSCAPES and PASCAL VOC in Section 5. Furthermore, we performed a number of additional studies to gain further insights into the effects of complementary labels, such as the dependence on the number of labels and the relative importance of label types. Our contributions are summarised as follows:

- We show empirical evidence that *inexpensive complementary labels improve model performance in low-data regimes*.
- We observe this effect in 1) a high-impact medical task where training examples are difficult to acquire and expert annotations are expensive and, 2) two well-studied public datasets.
- We conduct a series of studies that reveal further insights about this phenomenon. We show *a)* how the effect lessens as data increases *b)* that complementary labels provide robustness to annotator bias *c)* the effectiveness of different labels *d)* trivial labels are not useful *e)* performance increases with more labels *f)* low-quality labels are nearly as good as high-quality

labels *g)* complementary labels increase training stability *f)* complementary labels provide some robustness to domain shifts.

- We release the CSAW-S dataset used in this study to the public, which contains valuable mammography images with labels from multiple experts and non-experts that can be used to replicate our study and for other segmentation tasks.

Finally, to promote transparency and reproducibility, we share our open-source code, available at github.com/ChrisMats/seemingly_uninformative_labels and CSAW-S at <https://github.com/ChrisMats/CSAW-S>.

2. Related Work

Perhaps the most well-established method of dealing with insufficient training data is to learn transferable representations in a similar domain where data is more abundant, a technique routinely executed by means of pretraining on IMAGENET. The underlying assumption for this approach is that the domain gap is small, *i.e.*, the distribution the model is pretrained on has structural similarities to the target task’s conditional probability distribution (Bengio, 2012). Unfortunately, this assumption does not necessarily hold for all tasks (Azizpour et al., 2016), for example between natural and medical images. Raghu et al. have shown that IMAGENET pretraining offers only marginal improvement for some medical tasks, mainly attributed to better weight scaling and initialization (Raghu et al., 2019).

He et al. showed that initialization with IMAGENET yields no gains compared to random initialization in the big-data regime (He et al., 2019b). Interestingly, as they moved towards the low-data regime they noticed benefits from IMAGENET pretraining began to appear for PASCAL VOC but not for MS-COCO. The authors link this effect to the relatively low number of classes and object instances in the PASCAL VOC dataset. In other words, IMAGENET pretraining helps more when the downstream task is less diverse. Azizpour et al. argue that IMAGENET pretrained models benefit more from the diversity of classes than the number of training examples. Finally, it has been established that label correlation can increase accuracy in multi-label classification tasks by providing information regarding the interactions among them (Huang & Zhou, 2012). In line with these insights, we conjecture that more diverse annotations yield better representations and consequently better performance in scarce data settings.

Beyond transfer learning, modern methods such as unsupervised or weakly/semi-supervised approaches, can learn representations yielding comparable performance to supervised learning (Jing & Tian, 2019) provided the domain gap is relatively small. Recently, He et al. reported downstream

task performances even better than supervised pretraining (He et al., 2019a). Unfortunately, as in the supervised case, these methods require large amounts of data, which is a challenge in medical tasks (Hesamian et al., 2019).

Data augmentation techniques are widely employed when training models with limited data. In their seminal work, Ronneberger et al. have shown the effectiveness of suitable augmentations for medical image segmentation. Learning augmentations during training (Cubuk et al., 2019; Zhao et al., 2019) and GAN-based augmentation (Frid-Adar et al., 2018; Mondal et al., 2018) can also be employed in order to alleviate the effects of data scarcity.

Promising results have also been shown by k -shot methods (Roy et al., 2020) in low data regimes where certain classes only have a few or no representation in the training set. However, the domain gap and data scarcity are also significant problems in the k -shot settings as most of the methods in the literature rely on IMAGENET pretraining, which limits its applicability in non-natural image domains.

In this work, we argue that adding new labels that complement the ones provided for the principal task improves generalization in low-data regimes. Side information is a term for extraneous information, often a different modality than the source data, that can be exploited for a principal task (Kang et al., 2017). For example, (Tian et al., 2015) used additional datasets to introduce side information—expressed as scene and pedestrian attributes—to pedestrian classification.

3. Complementary Labels

Our central idea is simple but effective in practice. We argue that seemingly uninformative complementary labels, used as additional learning targets, have a direct impact on the model’s generalization for image segmentation in low data regimes.

These new annotations are often *seemingly uninformative* in the sense that they do not contain any direct information about the object of interest. However, they do contain useful information describing other semantically meaningful objects present in the image. These labels complement the expert labels by providing rich contextual information, so we refer to them as *complementary labels*. For the task of locating tumors in mammograms, complementary labels might include the *skin*, *pectoral muscle*, and other parts of anatomy (Figure 2). Complementary labels do not require expert knowledge and do not need to be particularly accurate, so crowdsourcing or other low-cost solutions can be employed to collect them.

The advantage of complementary labels is that they are easy to obtain and the performance boost comes without major changes to the model architecture or additional training ex-



Figure 2. The CSAW-S dataset, released here, contains 342 mammograms with expert radiologist labels for cancer and *complementary labels* of breast anatomy made by non-experts. The non-expert labels are imperfect and in some cases may seem uninformative, yet they provide useful cues for segmentation of the tumor.

amples. On the other hand, as more expert data are available, complementary labels yield diminishing returns (Figure 1). So consideration must be given to the performance-cost trade-off between obtaining inexpensive labels for existing data or collecting new data with expert labels.

Our observation invites an obvious question – *why* does adding complementary annotations lead to better performance? Although we cannot provide a definitive answer, we offer two plausible explanations:

1. *Complementary labels encourage learning of enriched representations.* With sufficient training examples, deep neural networks learn rich features that represent not only the object of interest, but also model the diversity of shapes and textures in the background. In classification and detection tasks, evidence has shown that networks learn to exploit information from the background (Ribeiro et al., 2016; Ghorbani et al., 2019). But in the low data regime, the network struggles to model the background because the data is insufficient to capture such a diverse distribution. Complementary labels help the network make sense of reduced background data by structuring it into more meaningful sub-classes with less individual variation. Analogous explanations have been hypothesized in (Huh et al., 2016; Azizpour et al., 2016) for classification and object detection (He et al., 2019b), where it is argued that providing fewer labels per image has a similar effect to removing training examples.
2. *Complementary labels help to model interactions between objects.* Objects that are near, interact with, or look similar to the target object contain information regarding the correlation and interaction that the target has with its environment. Providing complementary labels allows the network to exploit interactions and correlations between these additional labels and the target, resulting in better generalization. A similar line of

thinking explains how knowledge distillation benefits from interactions between labels through soft labeling (Hinton et al., 2015; Furlanello et al., 2018). Likewise, in (Huang & Zhou, 2012) correlations between rare labels are exploited to achieve better performance in a multi-label prediction task.

In the following sections, we provide experimental evidence supporting the benefits of complementary labels in low data regimes and explore practical questions surrounding their use and effects on generalization. In particular, we address:

- What is the effect of changing the number of training examples/expert annotations?
- What constitutes a good complementary label?
- What is the effect of adding new complementary labels?
- How does the quality of complementary labels affect performance?
- Do complementary labels provide robustness: 1) to annotator bias, 2) to domain shifts in the training data, and 3) training stability?

Through the following experiments on CSAW-S and two well-known public datasets, we attempt to characterize these properties of complementary labels.

4. Experiments on Medical Images

In this section we investigate how the addition of complementary labels affects the model’s performance in a low data setting on an expert medical task.

4.1. The CSAW-S Dataset

The CSAW-S dataset is a companion subset of CSAW, a large cohort of mammography data gathered from the entire population of Stockholm invited for screening between 2008 and 2015, which is available for research (Dembrower et al., 2019). We release the CSAW-S subset containing mammography screenings from 172 different patients with annotations for semantic segmentation. The patients are split into a test set of 26 images from 23 patients and training/validation set containing 312 images from 150 patients. Further details regarding the collection and pre-processing of the data can be found in the Appendix.

The training/validation images are accompanied by cancer annotations by an expert radiologist, EXPERT 1, and the test images come with cancer annotations from two additional radiologists, EXPERT 2 and EXPERT 3. Complementary labels are provided for the entire dataset in the form of full pixel-wise masks of each image corresponding to 11 additional highly imbalanced classes representing breast

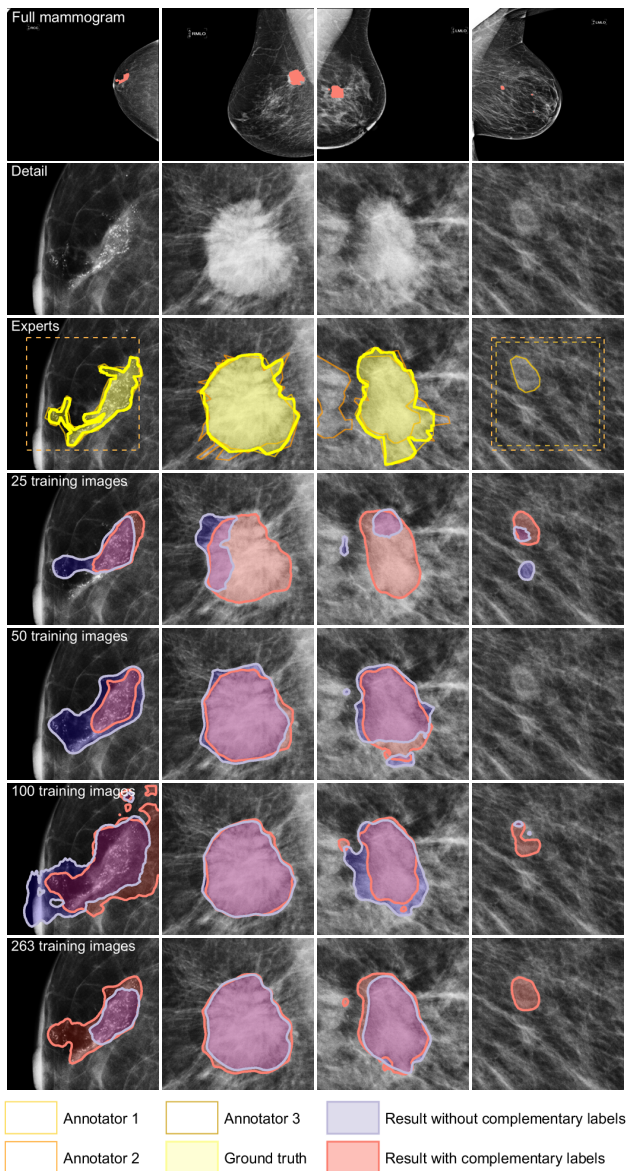


Figure 3. Segmentation results on CSAW-S show complementary labels improve performance in low-data settings. From top to bottom: the full mammogram, a detailed image of the tumor region, the expert annotations and ground truth (yellow), predictions from networks trained with only the tumor labels (blue) and predictions from a network trained with both tumor labels and complementary labels (red). Experiments are repeated using small training sets of varying size, $N = \{25, 50, 100, 263\}$, the best result from 5 runs is shown for each model. Expert annotations appear in shades of gold; a dashed box indicates an annotator did not find a tumor. Ground truth is determined where at least two annotators agree cancer is present (consensus in column 4 is *no tumor*).

anatomy and other objects (see Figure 2 and Appendix for details). *The complementary annotations were sourced from non-experts with no medical training, and therefore contain errors.* Complementary annotations from EXPERT 1 and EXPERT 2 are also provided for the test set. The ground

Table 1. Expert and model agreement on CSAW-S (IoU)

	EXPERT 1	EXPERT 2	EXPERT 3	with comp.*	without comp.*
Ground Truth	0.69	0.66	0.83	0.33	0.21
EXPERT 1		0.67	0.68	0.31	0.21
EXPERT 2			0.66	0.29	0.17
EXPERT 3				0.34	0.21

*Mean IoU for models trained with 50 patients.

truth for the expert task is determined as regions where at least two experts agree cancer is present (Fig. 3).

4.2. Experimental Setup

We split the train/validation sets by patient, 130/20. This resulted in 263/49 images per set. Many of our experiments investigate how performance depends on the size of the training set. For these cases, we create subsets of $N = \{25, 50, 100, 263\}$ images by sampling the training set at the patient level. We perform five runs for each experiment by sampling with replacement 5 times, and report the average performance along with the 95% confidence interval.

Our goal is to improve expert task performance by including complementary non-expert labels. Hence, we define two training settings to test this

1. *with complementary labels* – where expert labels and complementary labels are provided to train the network
2. *without complementary labels* – where only expert labels are provided.

These cases are compared throughout Sections 4 and 5.

4.3. Implementation Details

We use DeepLab3 (Chen et al., 2017) with ResNet50 (He et al., 2016) as the backbone for all experiments. Following He et al. and Raghu et al., we initialize all models with IMAGENET pretrained weights and we replace BatchNorm layers with GroupNorm layers (Wu & He, 2018). We use an ADAM (Kingma & Ba, 2014) optimizer throughout our experiments. Due to memory limitations and the high resolution of mammograms, we train using 512×512 patches. To ensure good representation in the training data, for every full image we sample a center-cropped patch from 10 random locations belonging to each of the 12 classes (the same for training with and without complementary labels).

To alleviate overfitting issues associated with extreme low data regimes, we employ an extensive set of augmentations including rotations and elastic transformation in addition to standard random flips, random crops of 448×448 , random brightness and random contrast augmentations. We report results for each run using the best checkpoint model. Since

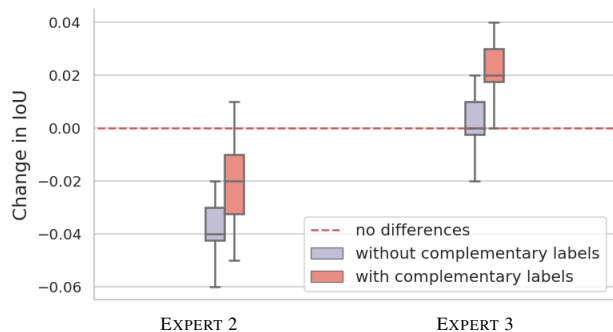


Figure 4. Complementary labels provide robustness to annotator bias. The training data was annotated by EXPERT 1, biasing the models towards this expert. Test results evaluated on annotations from the other experts show that models provided with complementary labels performed consistently better (indicating robustness to bias) than those without (measured by change in IoU).

the cross entropy loss does not precisely represent the IoU metric we consider both the validation IoU and loss when selecting the best model. For all of our experiments we fine-tuned the learning rate for each setting and the results are averaged over 5 runs, unless otherwise specified.

4.4. Results

Do complementary labels help? Our main results appear in Figure 1, where we quantify the effects of adding complementary labels measured by IoU. Evidently, training with complementary labels outperforms the case where only expert tumor annotations are used by a large margin. Although the absolute IoU performance is relatively low, inter-expert agreement is also low (≈ 0.67), and adding complementary labels results in a 57% relative gain in IoU towards the level of expert agreement (Table 1). Segmentation visualizations appear in Figure 3, showing especially noticeable benefits when very little training data is available.

What is the effect of changing the amount of training examples? From the results in Figure 1 it is evident that the benefit of complementary labels is magnified as we move towards lower data regimes. Models trained with expert and complementary labels from 50 patients outperform models trained on only expert labels from 130 patients.

Do complementary labels provide robustness against annotator bias? A major challenge facing intelligent diagnostic systems is disagreement between experts (Kerlikowske et al., 1998). Because expert annotations are costly, most medical imaging datasets are annotated by only a few experts. The CSAW-S training set was annotated solely by EXPERT 1. This can introduce bias in the model towards the opinion of that expert. To test if adding complementary labels introduces robustness to annotator bias in our models, we measure performance using test set annotations from the other experts. In Figure 4 we report the change in IoU when

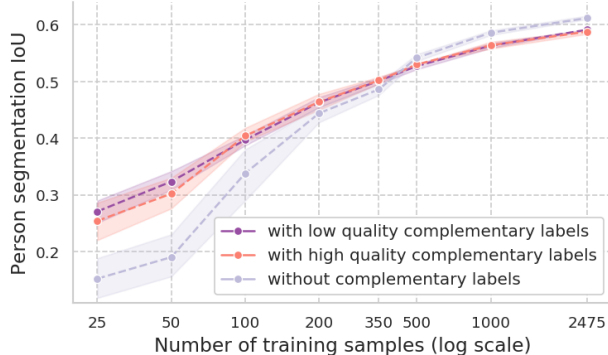


Figure 5. Results on CITYSCAPES using the *person* class as the expert class. Performance with (red) and without (blue) complementary labels is measured for varying numbers of training images over five repetitions (shaded area indicates 95% CI). An additional comparison to examine the impact of low-quality complimentary labels (violet) shows little effect.

evaluating using annotations from EXPERT 2 and EXPERT 3 compared to EXPERT 1. In both cases, adding complementary labels increases robustness against annotator bias. Interestingly, *both* models performed better when evaluated on EXPERT 3’s annotations.

5. Further Studies Using Public Data

We confirm our findings on two publicly available datasets, CITYSCAPES and PASCAL VOC, and also delve deeper – using these well-known datasets to investigate several interesting properties of complementary labels.

5.1. Public Datasets

CITYSCAPES (Cordts et al., 2016), is a sizable public dataset of urban street scenes with a large and diverse set of annotations. While the protocol for most segmentation research is to use 19 of the 34 annotated classes, in the interest of better understanding complementary labels, we use all 34 classes.

Although expert knowledge was not necessary to annotate CITYSCAPES, our goal is to test the effects of complementary labels on an expert task. We simulate this by choosing the *person* class as the target class, so the goal becomes person segmentation. The other 33 classes are either treated as complementary labels, or merged into a single *background* class for the baseline. We selected *person* because of the high intra-class variance, relatively small size, and moderate image frequency. Additionally, the *rider* class (e.g. cyclists, motorcycle riders) presents an opportunity to examine performance when confusing objects are present. We randomly select 500 images from the official training set to use as the validation set, and we use the rest for training. For testing, we use the official validation set of fine annotations.

PASCAL VOC (Everingham et al., 2015) is a well-studied

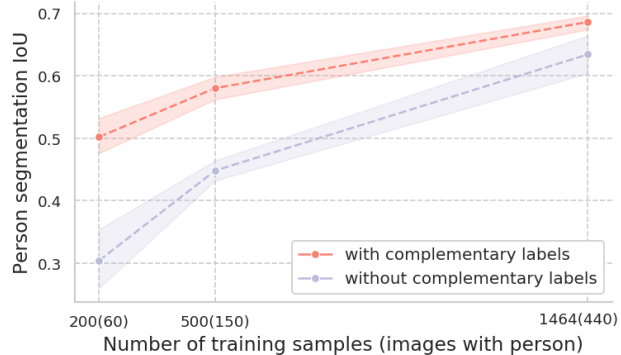


Figure 6. Results on PASCAL VOC using the *person* class as the expert class. Because *person* occurs rarely, measurements are taken at different numbers of training samples than the previous experiments.

object recognition dataset composed of 21 object classes and a background class. Here, we also choose the *person* as the main target and we consider the other 20 classes as either complementary labels or background. In PASCAL VOC, the *person* class is present in only 30% of the training images, which allows us to investigate how complementary labels perform for expert tasks on rare events. We used the official `train-2012` as training set and we sampled 500 images from the `test-2012` as our validation set. We used the remaining 949 images from `test-2012` to evaluate the final performance.

5.2. Implementation Details

For consistency, we kept the training procedure constant throughout our experiments and used the training settings outlined in Section 4, with the following exceptions: random rotation and elastic augmentations were omitted, and the full images were resized to 512×512 instead of using patches. We note that settings optimized for CSAW-S might be sub-optimal for these datasets (see Appendix), but we opted to maintain the same settings for better comparability.

5.3. Results

Are the benefits of complementary labels observed in other data? We confirm our previous findings on medical images with data from CITYSCAPES and PASCAL VOC in Figures 5 and 6 respectively. We note that, although they are not directly comparable, the IoU gap between models with/without complementary labels is nearly the same for CSAW-S and CITYSCAPES, and the gap nearly doubles for PASCAL VOC. The size of these datasets allowed us to explore how the trend evolves further from the low-data regime. Interestingly, on CITYSCAPES there appears to be a crossover when the network is provided with approximately 400 examples where complementary labels seem to hurt performance rather than help.

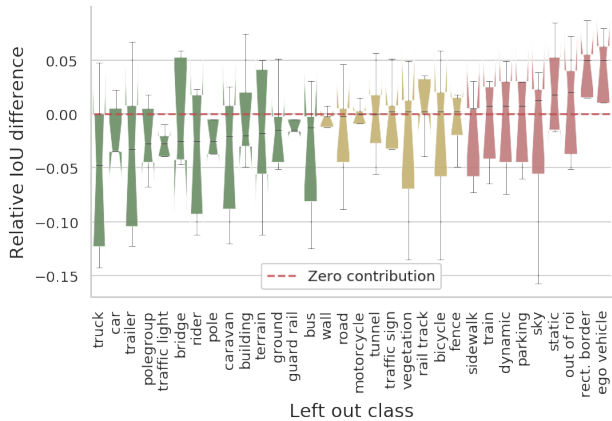


Figure 7. Contributions of individual complementary labels. We conducted a leave-one-out experiment where the network was trained in turn using every complementary label in CITYSCAPES except one. The importance of each label is measured by the relative change in IoU for person classification when the label is omitted. Labels marked in green, which tend to occur frequently in the data with diverse appearances, improve performance. Labels marked in yellow have little discernible effect, and labels marked in red hurt performance. These tend to be omnipresent and trivial (e.g. ego vehicle and rectification border).

Do all complementary labels contribute equally? We investigate whether all complementary labels contribute positively, and to what degree each complementary label helps. To this end, we conducted a leave-one-out experiment where the network was trained with every complementary label except one, which was merged to the background. In this way, we can measure the importance of the left-out class by the change in segmentation IoU. As seen in Figure 7, there is a clear disparity between various complementary labels. Most of the complementary labels contribute positively – removing them hurts the network’s performance (green). The effect is unclear for nine (yellow), and for the rest there is a clear advantage in removing them (red).

Looking more closely at Figure 7, we infer that the labels which seem to contribute most are those that appear frequently with diverse appearances (truck, car, traffic light). Also, smaller objects that could potentially be confused with the *person* class appear to help (pole groups). The least helpful labels seem to be omnipresent and trivial (e.g. ego vehicle and rectification border). Finally, we note that *rider*, which can be easily confused with *person*, is important for the network to avoid false positives (further experiments in Appendix).

What is the effect of adding new complementary labels? To measure the impact that the number of complementary classes has on the model’s performance, we conducted a series of experiments on CITYSCAPES where the number of complementary labels provided to the network is steadily

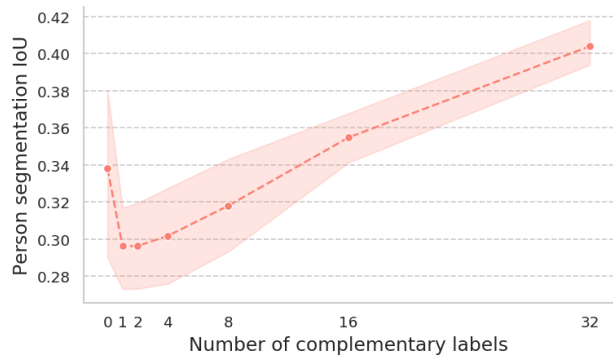


Figure 8. The effect of adding new complementary labels. We measure how IoU performance changes on CITYSCAPES for a steadily increasing number of randomly selected complementary classes. Interestingly, adding only a few classes hurts performance, but once a sufficient number of labels is reached there is a clear advantage. These results, with those of Figure 7, suggest that the choice of which labels to include is important.

increased, $N = \{0, 1, 2, 4, 8, 16, 32\}$. The labels added in each run are randomly selected, and each experiment is repeated 5 times. The number of training examples was fixed to 100 and *person* was again used as the expert target. We observe an interesting trend in Figure 8, where adding only a few random complementary labels hurts performance. But once a sufficient number of labels are used, there is a clear advantage to adding them. These results, combined with those of Figure 7, imply that care should be taken when choosing which complementary labels to include. Either that, or a sufficient number should be collected to overcome adverse effects from certain classes.

How does the quality of complementary labels affect performance? We demonstrated in Section 3 that complementary labels do not require expert knowledge, and also noted that some the labels in CSAW-S contain errors and noise. This leads us to the question: does quality of the complementary annotations affect performance? Luckily, the CITYSCAPES dataset includes two sets of annotations for semantic segmentation, the fine and coarse set. The fine set (high quality), seen in Figure 9, consists of high quality dense pixel annotations whereas the coarse set (low quality) includes approximate polygonal annotations – which omit regions of objects in many cases and in some cases actually mislabel objects. We tested the performance of the model trained with low-quality annotations and compared it to the high-quality annotations. To ensure the expert task was not affected, we used high-quality annotations for the *person* class in both cases. As we can see from Figure 5, the complementary labels do not need to be accurate. Unexpectedly, for the extreme low regimes, the low-quality annotations resulted in slightly higher IoU scores. A possible explanation is that the coarse labels avoid boundary regions between objects, which may help learning.

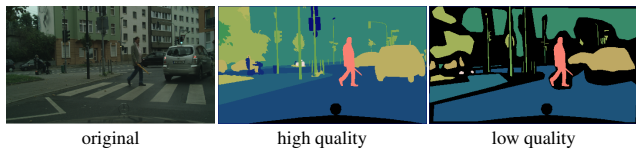


Figure 9. Testing robustness to label quality. To test how complementary label quality affects performance, we compare how models trained with fine labels fare against models trained with coarse labels on CITYSCAPES. Note that high quality labels are given for the expert class, *person*, in both cases. Results in Figure 5 show little effect from reducing quality of complementary labels.

Do complementary labels improve training stability?

Over the course of our experiments, we noticed that the addition of complementary labels resulted in more stable training in the low data regime. Training with complementary labels yields smoother and clearer learning curves, as seen in Figure 10, which makes it easier to identify signs of overfitting.

Do complementary labels provide robustness to domain shifts in the training data?

A common issue in many datasets – especially in medical applications – are domain shifts within the data. For example, medical images can be acquired from a small number of clinics with different devices. The CITYSCAPES images were collected from 18 different German cities and Zurich. This phenomenon can cause generalization issues if the training data is not representative of the true distribution.

As a final investigation, we test if adding complementary labels improves model robustness to domain shifts. We set up an experiment in which domain shifts are artificially imposed in the training data as follows. An ordered training set is created by shuffling images individually from each city, and then placing them in a random order grouped by city. For example, this may result in a training set with images from Stuttgart, then Aachen, Hamburg, etc. Then, we repeat the experiments for *person* segmentation. The models are trained using subsets of the randomly ordered and shuffled training sets, with the same schedule as our main experiments $N = \{25, 50, 100, 200, 350, 500, 1000, 2475\}$. We repeat this 5 times for each N . The result of this procedure is that models trained with $N = 25$ or $N = 50$ will only see data from a single city, but will be tested on a set containing all cities (each city contains between 77 and 259 images). In this way, we artificially impose a domain shift which lessens as more training data is added (more cities will appear).

Our results appear in Figure 5 and Figure 6 the Appendix. We find that adding complementary labels still improves performance in the presence of domain shifts. Although the absolute IoU performance is lower than in in Figure 5, the performance gap between models with and without complementary labels widens to some degree. Note that the curve for the model trained without complementary labels

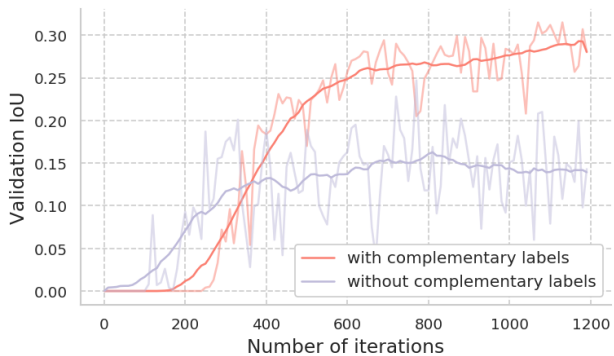


Figure 10. Complementary labels improve training stability in low data regimes. As seen in this IoU evolution curve training on 25 images from CITYSCAPES, training with complementary labels yields smoother and clearer learning curves, which helps to identify signs of overfitting.

shows more variance, and is more sensitive to steps where new cities are added.

6. Conclusions

In this work we consider a semantic segmentation task and show that adding inexpensive and seemingly uninformative labels can significantly increase the model’s generalisation under low data regimes. We demonstrate the effects of these complementary labels on an expert medical task and on the CITYSCAPES and PASCAL VOC datasets. We release a new dataset, CSAW-S, along with this study which contains valuable mammography images with labels from multiple experts and non-experts that can be used to replicate our study and for other tasks in conjunction with its larger companion dataset.

We identify several interesting properties of complementary labels. First, these labels yield larger benefits when data is scarce. This property is critical in domains such as medicine where data gathering and annotation costs are often prohibitive. Separately, we find that complementary labels need not be of high quality, which suggests crowd-sourcing solutions or automation may be utilized for additional cost savings. We note that not all labels are equally useful, a fact that can help guide the annotation design process – we witnessed that certain trivial objects hurt performance. Complementary labels seem to provide several forms of robustness to some degree: against annotator bias, against domain shift, and increased training stability.

Complementary labels are a simple and effective means to improve performance in low-data settings. We believe they should be an essential tool in every practitioner’s toolkit. However, further research is required to gain a better understanding of the mechanisms that result in this behavior.

Acknowledgements

This work was partially supported by the Wallenberg Autonomous Systems Program (WASP), the Swedish Research Council (VR) 2017-04609, and Region Stockholm HMT 20170802. We would like to thank Josephine Sullivan for the fruitful discussions. Finally, we thank the reviewers for their constructive criticism and comments.

References

- Azizpour, H., Razavian, A. S., Sullivan, J., Maki, A., and Carlsson, S. Factors of transferability for a generic convnet representation. *IEEE transactions on pattern analysis and machine intelligence*, 38(9):1790–1802, 2016.
- Bengio, Y. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pp. 17–36, 2012.
- Chen, L.-C., Papandreou, G., Schroff, F., and Adam, H. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., and Schiele, B. The cityscapes dataset for semantic urban scene understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 113–123, 2019.
- Dembrower, K., Lindholm, P., and Strand, F. A multi-million mammography image dataset and population-based screening cohort for the training and evaluation of deep neural networks—the cohort of screen-aged women (csaw). *Journal of digital imaging*, pp. 1–6, 2019.
- Everingham, M., Eslami, S. M. A., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, January 2015.
- Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J., and Greenspan, H. Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification. *Neurocomputing*, 321:321–331, 2018.
- Furlanello, T., Lipton, Z., Tschannen, M., Itti, L., and Anandkumar, A. Born again neural networks. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1607–1616, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- Ghorbani, A., Wexler, J., Zou, J. Y., and Kim, B. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, pp. 9273–9282, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019a.
- He, K., Girshick, R., and Dollar, P. Rethinking imagenet pre-training. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019b.
- Hesamian, M. H., Jia, W., He, X., and Kennedy, P. Deep learning techniques for medical image segmentation: Achievements and challenges. *Journal of Digital Imaging*, 32(4):582–596, Aug 2019.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- Huang, S.-J. and Zhou, Z.-H. Multi-label learning by exploiting label correlations locally. In *Twenty-sixth AAAI conference on artificial intelligence*, 2012.
- Huh, M., Agrawal, P., and Efros, A. A. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.
- Jing, L. and Tian, Y. Self-supervised visual feature learning with deep neural networks: A survey. *arXiv preprint arXiv:1902.06162*, 2019.
- Kang, D., Dhar, D., and Chan, A. Incorporating side information by adaptive convolution. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 3867–3877. Curran Associates, Inc., 2017.
- Kerlikowske, K., Grady, D., Barclay, J., Ernster, V., Frankel, S. D., Ominsky, S. H., and Sickles, E. A. Variability and accuracy in mammographic interpretation using the american college of radiology breast imaging reporting and data system. *Journal of the National Cancer Institute*, 90(23):1801–1809, 1998.

- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Mondal, A. K., Dolz, J., and Desrosiers, C. Few-shot 3d multi-modal medical image segmentation using generative adversarial learning. *arXiv preprint arXiv:1810.12241*, 2018.
- Raghu, M., Zhang, C., Kleinberg, J., and Bengio, S. Transfusion: Understanding transfer learning for medical imaging. In *Advances in Neural Information Processing Systems*, pp. 3342–3352, 2019.
- Ribeiro, M. T., Singh, S., and Guestrin, C. ” why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- Ronneberger, O., Fischer, P., and Brox, T. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Roy, A. G., Siddiqui, S., Pölsterl, S., Navab, N., and Wachinger, C. ‘squeeze & excite’ guided few-shot segmentation of volumetric images. *Medical image analysis*, 59:101587, 2020.
- Tian, Y., Luo, P., Wang, X., and Tang, X. Pedestrian detection aided by deep learning semantic tasks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- Wu, Y. and He, K. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19, 2018.
- Zhao, A., Balakrishnan, G., Durand, F., Guttag, J. V., and Dalca, A. V. Data augmentation using learned transformations for one-shot medical image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8543–8553, 2019.