

# Learning with Weak Supervision for Email Intent Detection

Kai Shu\*  
Arizona State University  
Tempe, AZ  
kai.shu@asu.edu

Subhabrata Mukherjee†  
Microsoft Research AI  
Redmond, WA  
submukhe@microsoft.com

Guoqing Zheng†  
Microsoft Research AI  
Redmond, WA  
zheng@microsoft.com

Ahmed Hassan Awadallah  
Microsoft Research AI  
Redmond, WA  
hassanam@microsoft.com

Milad Shokouhi  
Microsoft  
Bellevue, WA  
milads@microsoft.com

Susan Dumais  
Microsoft Research AI  
Redmond, WA  
sdumais@microsoft.com

## ABSTRACT

Email remains one of the most frequently used means of online communication. People spend significant amount of time every day on emails to exchange information, manage tasks and schedule events. Previous work has studied different ways for improving email productivity by prioritizing emails, suggesting automatic replies or identifying intents to recommend appropriate actions. The problem has been mostly posed as a supervised learning problem where models of different complexities were proposed to classify an email message into a predefined taxonomy of intents or classes. The need for labeled data has always been one of the largest bottlenecks in training supervised models. This is especially the case for many real-world tasks, such as email intent classification, where large scale annotated examples are either hard to acquire or unavailable due to privacy or data access constraints. Email users often take actions in response to intents expressed in an email (e.g., setting up a meeting in response to an email with a scheduling request). Such actions can be inferred from user interaction logs. In this paper, we propose to leverage user actions as a source of weak supervision, in addition to a limited set of annotated examples, to detect intents in emails. We develop an end-to-end robust deep neural network model for email intent identification that leverages both clean annotated data and noisy weak supervision along with a self-paced learning mechanism. Extensive experiments on three different intent detection tasks show that our approach can effectively leverage the weakly supervised data to improve intent detection in emails.

## ACM Reference Format:

Kai Shu, Subhabrata Mukherjee, Guoqing Zheng, Ahmed Hassan Awadallah, Milad Shokouhi, and Susan Dumais. 2020. Learning with Weak Supervision for Email Intent Detection. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*,

\*Research was conducted while interning at Microsoft Research AI.

†Equal contributions.

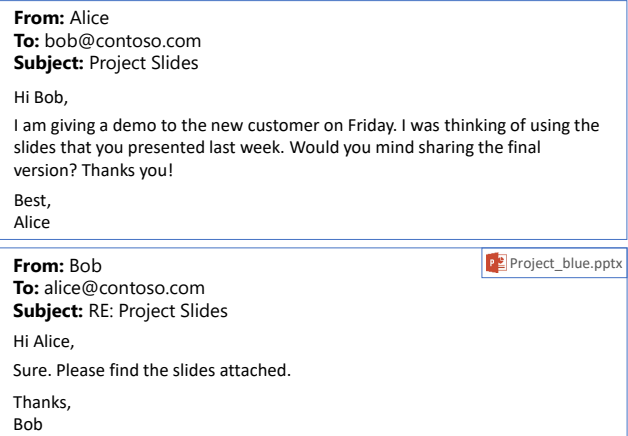
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

*SIGIR '20, July 25–30, 2020, Virtual Event, China*

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-8016-4/20/07...\$15.00

<https://doi.org/10.1145/3397271.3401121>



**Figure 1: An illustration of user interaction for email intent classification. On the upper side, the sender Alice is requesting information from the recipient Bob. On the lower side, the recipient Bob is replying Alice with an attachment – an interaction that can be leveraged for weak supervision.**

July 25–30, 2020, Virtual Event, China. ACM, New York, NY, USA, 10 pages.  
<https://doi.org/10.1145/3397271.3401121>

## 1 INTRODUCTION

Email has continued to be a major tool for communication and collaboration over the past decades. The volume of email messages exchanged daily has also continued to grow and is projected to reach 306.4 billion messages and 319.6 billion messages a day by the end of 2020 and 2021 respectively [1]. In addition to the significant volume of messages, email is one of the top time consuming activities for information workers. Recent studies show that communicating with colleagues and customers takes up to 28% of information workers' time, second only to role-specific tasks at 39% [9].

Such widespread use and significant amount of time spent on email have motivated researchers to study how people use email and how intelligent experiences could assist them to be more productive [12, 24, 28]. One of the earliest works to characterize the main purpose email serves in work settings is that of Dabbish et al. [12]. They conducted a survey of 124 participants to characterize different aspects of email usage. Based on this, they identified four distinct uses of email: task management, social communication, scheduling, and information exchange. More recent work [50]

conducted a large scale analysis of enterprise email identifying several use cases with many sub intents such as requesting an action, promising an action, updating a meeting, requesting information, social interaction, etc. Many other studies have focused on proposing methods for detecting intent of or suggesting actions in response to an email [6, 11, 28, 46]. Detecting intents in communications can integrate machine intelligence into email systems to build smart email clients that provide more value to email users. Several such applications have been studied including creating intelligent experiences that offer to assist users with scheduling a meeting [20], detecting action items [6], automatically populating to-do lists [29], creating alerts for high-priority messages [21], sharing documents [51], and answering questions [53].

Previous work posed email intent classification as a supervised learning problem where human annotators were asked to annotate email messages given a predefined taxonomy and machine learning models were built to identify intents using the annotated dataset for training. Supervised models, especially those employing deep neural networks, rely on large scale annotated training data for learning. In many applications, manual annotation is either expensive and time-consuming to acquire, or infeasible due to privacy concerns for sensitive data. This is exactly the case for email data since its personal and private nature makes it hard to collect human annotations for. Even when annotations are collected, they are done on a limited amount of data that may not be sufficient or representative of the different domains.

Many application domains like recommendation, search, and email communication have rich user activities and interactions that can provide additional signals for learning [2, 3, 22, 26]. For example, leveraging user interaction (e.g., clicks) for web search ranking has been extensively studied [2]. Most email clients also allow users to manage their calendars, task lists, etc. Users interact with these information in different ways including responding, forwarding, flagging emails, setting up appointments, etc. Many of these user actions are directly correlated with the intent of the email. For example, many scheduling intents could be correlated with taking an action on the user’s calendar such as creating or updating a calendar item. These actions may correlate to a certain intent but are also noisy in nature. Refer to Figure 1 for an example. Consider the scenario where we want to detect emails where the sender is requesting a document from the recipient. In the absence of enough annotated examples, we may heuristically label all emails that received a response email *with attachments* as positive examples. However, there would be a lot of false positives since users may send attachments even without being requested for. Similarly, there will also be many false negatives, as users may not always send attachments even when requested to do so. In this work, we pose the email intent detection as a weakly supervised learning problem. We assume, as is usually the case in practice, that we have access to a small amount of annotated training examples for any given intent and a lot of weakly labeled noisy instances constructed from interaction-based heuristics (such as receiving a reply with an attachment). Note that such interaction signals typically cannot be used as features since they are only available after the user has processed and interacted with the message (e.g., predicting that one should schedule a meeting in response to a message after she has already done so is useless). Also, note that these weak supervision

signals can be used by a model for training without requiring any human inspection alleviating a lot of privacy concerns. We propose an end-to-end trainable framework with deep neural networks to model these two sources of supervision simultaneously.

In summary, this paper makes the following contributions:

- **Application.** We show that weak supervision from user interaction is effective in the presence of limited amount of annotated data for the task of email intent identification. This pushes the frontier on weak supervision on email-related tasks where the focus has traditionally been on training fully supervised models.
- **Model.** We propose a unified framework Hydra to leverage cleanly annotated examples and weakly labeled ones *jointly* by embedding them in a shared representation space. We incorporate ideas from prior work on label correction to better handle noisy labels.
- **Learning.** We propose a learning mechanism for Hydra based on prior works in curriculum learning and self-paced learning [27] to judiciously select informative weak instances to learn from.
- **Experiments.** We conduct extensive experiments on real-world datasets to demonstrate the effectiveness of the proposed approach – obtaining an accuracy improvement of 3% to 12% on average over state-of-the-art methods for different settings.

## 2 EMAIL INTENT DETECTION

In this section, we first formulate the weakly supervised learning problem of intent detection for email communication. Thereafter we describe in details the task, dataset, and how weak supervision can be leveraged from user interactions to help in the task.

### 2.1 Problem Statement

Let  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$  denote a set of  $n$  email messages with manually annotated clean labels, with  $\mathcal{X} = \{x_i\}_{i=1}^n$  denoting the messages and  $\mathcal{Y} = \{y_i\}_{i=1}^n$  the corresponding clean labels. Each message  $x_i = \{w_1^i, \dots, w_{m_i}^i\}$  contains a sequence of  $m_i$  words. In addition to the small set of labeled examples, there is a large set of unlabeled examples. Usually the size of the clean labeled set  $n$  is much smaller than the unlabeled set due to labeling costs or privacy concerns for email data. For the widely available unlabeled samples, weak labels can be obtained based on *user interactions* with emails (as illustrated in the example for Figure 1 and more details in Section 2.3). Denote the weakly labeled set by  $\tilde{\mathcal{D}} = \{\tilde{x}_j, \tilde{y}_j\}_{j=1}^N$  where  $\tilde{\mathcal{X}} = \{\tilde{x}_j\}_{j=1}^N$  denotes the set of  $N$  unlabeled messages and  $\tilde{\mathcal{Y}} = \{\tilde{y}_j\}_{j=1}^N$  denotes the set of weak labels derived from user interactions. We now formally define our problem as:

**Problem Statement:** Given a small manually annotated data  $\mathcal{D}$  and a large set of weakly labeled data  $\tilde{\mathcal{D}}$  with weak labels derived from user interactions, learn an intent classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}$  which generalizes well onto unseen samples.

### 2.2 Dataset

To better understand how real-world enterprise email data exhibits user intents, we leverage the Avocado<sup>1</sup> dataset [38], which contains an anonymized version of the Outlook mailbox for 279 employees

<sup>1</sup>Avocado is a more appropriate test bed than Enron [25] since it contains additional meta-data and information beyond email such as calendar, tasks, etc.

with various meta information. The full Avocado corpus contains 938, 035 emails, 26, 980 meeting schedules, and 325, 506 attachments. We focus on multiple *intent detection* tasks on this data and accordingly devise weak labeling functions from user interactions.

**Intent types:** Prior works [13] have categorized email intents into four major categories: information exchange, task management, scheduling and planning, and social communications. Each category can have multiple fine-grained intents [51]. For instance, in the case of information exchange, *requesting information* is a common intent that indicates the sender is requesting information that can be potentially responded to by sharing a document. *Schedule meeting* refers to the sender’s intention to organize an event such as a physical meeting or a phone call, which belongs to the broader intent of scheduling and planning. In the case of task management intent, *promise action* is an intent that indicates the sender is committing to complete a future action. In this work, we focus on these three intents – request information, schedule meeting, and promise action (denoted by RI, SM, and PA, respectively). For instance, identifying the intent of requesting information allows an intelligent assistant system to automatically suggest files to share with the requester. This can result in improving the overall user experience and also user productivity. Table 1 shows some examples.

## 2.3 Deriving Weak Labels from User Interactions

With human annotations being hard to collect, if not impossible to obtain, in large scale, it is often cheaper and more beneficial to leverage weak supervision to build supervised models, particularly those based on deep neural networks. For email data, such weak supervision can be derived from user interactions. In this subsection, we discuss details to automatically obtain such weak labels from user interactions by using labeling functions and performing human evaluations to assess the quality of the weak labels.

**2.3.1 Weak Labeling Functions from User Interactions.** For each of the aforementioned intent types, we define the weak labeling functions as follows:

*Request Information (RI):* We observe that the action of replying with attachments may potentially indicate the email it replies to has the intent of RI. For example, the email “*Please forward me the final version for the slides*” is asking the recipient(s) to send a file back to the sender. Now, if a user replies with an email “*Please find the paper draft as attached*” along with an attachment, then the replied-to email is likely to contain the RI intent. However, this rule will have false positives since a user may reply with attachments even without being asked for. Additionally, messages with an RI intent may not receive a reply with an attachment or even any reply. Formally, the *weak* labeling function is:

**reply\_with\_attachment:** *If an email a is replying to another email b with an attachment, then email b is weakly-labeled with the RI intent.*

Note that we ignore the trivial attachments that are not likely to contain information related to RI intent (e.g. contact information, signatures, images, etc.).

*Schedule Meeting (SM):* Since we have access to not only user emails but also their calendars, we explore the temporal footprints of the scheduled meetings including the subject line of the meeting,

time, location and attendees. However, the emails that propose meeting requests are not directly associated with the schedule information. Therefore, we take the subject lines of the schedules as a query and search the emails that contain similar subject lines. This reveals the confirmation emails sent after someone accepted the meeting request. We temporally order the sent email together with the confirmation, and treat the earlier one as having the SM intent. The corresponding weak labeling function is defined as:

**confirmed\_schedule:** *If an email a has the same subject line with another email b confirming the schedule where a precedes b in the timeline, then a is weakly-labeled with the SM intent.*

*Promise Action (PA):* Outlook allows users to maintain a task list of items they need to do later. Tasks can be added to the task list by either directly creating them or by *flagging* emails that may contain future action items. The flags could be added by either the sender or the recipient of an email. We use this behavior as a proxy label for future action’s intent. For example, given an email from the sender as “*Would you be able to present your work in the meeting next week?*” with the urgency flag set and a response email as “*I can do this next week*” – we consider the latter email to have the PA intent. The corresponding weak labeling function is defined as:

**urgency\_reply:** *If an email a has replied to an email b which had a follow-up flag set, then a is weakly-labeled with the PA intent.*

It is worth mentioning that although these labeling functions are geared for intent detection for email communication, it is possible to devise similar or more sophisticated rules for different intents and even different domains.

**2.3.2 Quality of Weak Labeling Functions.** We conduct a manual evaluation to measure the quality of the labels generated by the weak labeling functions. We apply the aforementioned labeling functions (see Figure 1) to the Avocado corpus, and obtain the following weakly labeled *positive* instances: 8,100 emails for RI, 4,088 emails for SM, and 2,135 emails for PA. In addition, we treat the emails discarded by the weak labeling functions as *negative* instances. For each intent, we sample the same amount of negative instances as the positive ones to construct a balanced dataset. Note that the total amount of weak labeled data depends on the overall size of the email collection, how prevalent an intent is and the trigger rate of the labeling functions. In practice, developers may have access to a much larger pool of unannotated emails compared to the Avocado dataset containing mailboxes for only 279 users. This may enable generating even larger weakly labeled instances and potentially further improving the overall performance. We study the effects of the relative size of the clean and the weak data later in the experiments section.

To assess the quality of these weakly-labeled instances, we randomly select 100 emails from each of the positive and negative weakly-labeled sets that are sent for manual annotation. Each email is annotated by three people and their majority consensus is adopted as the final annotation label for the instance. Table 2 shows the confusion matrix from manual annotation. The accuracy of the weak labeling functions for the three intents RI, SM, and PA are 0.675, 0.71, 0.63, respectively. We observe that the accuracy of the labeling functions, while far from perfect, is also significantly better

**Table 1: Examples of different intent types in enterprise emails with weak labeling rules derived from user interactions.**

Intent	Example	Weak Supervision Rule
Request Information (RI)	Please forward me the final version for the slides	reply_with_attachment
Schedule Meeting (SM)	Let’s meet on Monday to discuss these accounts	confirmed_schedule
Promise Action (PA)	Once we have made reservations, we will let you know	urgency_reply

**Table 2: Confusion matrix for human evaluation of weak labeling functions.**

Intent	Predictions	True Positive	True Negative
RI	Positive	36%	64%
	Negative	1%	99%
SM	Positive	46%	54%
	Negative	4%	96%
PA	Positive	31%	69%
	Negative	5%	95%

**Table 3: Email datasets with metadata and user interactions. Clean refers to manually annotated emails, whereas weak refers to the ones obtained leveraging user interactions.**

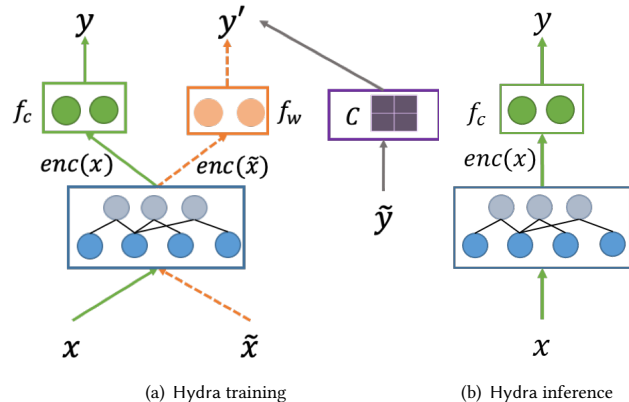
Email	Intent	Train		Dev	Test
		Clean	Weak		
Avocado	RI	1,800	16,200	334	336
	SM	908	8,176	1,008	1,010
	PA	474	4,270	518	518
Enron	SM	908	-	908	908

than random (0.5) for binary classification. This indicates that the weak labeling functions carry a useful signal albeit it being noisy.

## 2.4 Incorporating Small Amount of Clean Labels

Training neural networks with only noisy labels is challenging since they have high capacity to fit and memorize the noise [55]. Hence, it is useful to also incorporate clean labeled data in the training process. However, as discussed earlier, it is hard to collect human annotated data at scale especially when we need to support many intents across thousands of domains/organizations due to resource and privacy constraints. To better reflect these constraints in our experimental setup, we collected manual annotations for different intents such that the annotated set constitutes at most 10% of all the labeled samples (weakly as well as manually labeled). Note that the clean labeled and weakly labeled email instances have no overlap.

We first selected 10,000 Avocado email threads randomly and excluded them from the weak supervision data collection. Given the number of weakly labeled instances for each intent in Section 2.3, we collected manual annotations on 10% of the selected email threads. To this end, three annotators examined all the messages in each thread and annotated the first message with one or more of the intents as described above with majority votes deciding the final label. The Cohen’s kappa coefficient for inter-annotator agreement for each task was greater than or equal to 0.61 indicating a substantial agreement among the annotators [10]. The statistics



**Figure 2: Proposed framework Hydra for learning with weak supervision from user interactions. (a) Hydra leverages a label correction component (purple box) to rectify weak labels ( $\tilde{y} \rightarrow y'$ ) while learning; (b): During inference, Hydra uses learned feature representation (blue box) and function  $f_c$  to predict labels for (unseen) instances in test data.**

of the dataset used for our intent classification tasks are reported in Table 3. For each intent, there are more negative samples than positive ones. We down-sample the negative class to make the classes balanced for each task. Although the entire Avocado email collection is large, we have only a few manually annotated clean samples for each intent. This also motivates the idea of incorporating user interactions as weak supervision to build better predictive models than using the clean samples alone. Based on our dataset construction, clean labels account for 10% and weak labels constitute 90% of all the labeled samples (both clean and weak). Additionally, we report the performance of various methods in the extreme case – where clean labels account for only 1% of all the labeled samples obtained by further down-sampling the clean examples.<sup>2</sup>

## 3 JOINT LEARNING WITH CLEAN AND WEAK SOURCES OF SUPERVISION

Having defined the problem setting for intent detection with weakly supervised learning in the presence of a small set of cleanly labeled examples and a large set of weakly labeled ones, we now propose our approach to leverage these two sources of supervision jointly to learn an end-to-end model.

### 3.1 Hydra: Dual-source Supervised Learning

Multi-source learning has shown promising performance in various domains such as truth discovery [16], object detection [39], etc. In our email intent detection scenario, we have two distinct sources

<sup>2</sup>We will make the code and data publicly available at <https://aka.ms/HydraWS>, in accordance with the sharing policy for the Avocado dataset.

(dual-source) of supervision: clean labels coming from manual annotation and weak labels coming from heuristic labeling functions based on user interaction signals.

Our objective is to build a framework that leverages signals coming from both sources of supervision and learn an underlying common representation from the context. We develop a deep neural network where the lower layers of the network learn *common* feature representations of the input space (text of messages in our context), and the upper layers of the network *separately* model the mappings to each of the different sources of supervision. In this way, we are able to *jointly* leverage both the correlation and distinction between the clean and weak labels. Since the weak labels are obtained from labeling functions defined over user interactions, they contain complementary information to the clean labels annotated from message contents. We term this framework Hydra<sup>3</sup>.

Recall  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$  and  $\tilde{\mathcal{D}} = \{\tilde{x}_j, \tilde{y}_j\}_{j=1}^N$  to be the clean labeled data (based on manual annotation) and the weak labeled data (based on user interactions) respectively. Let  $enc(x; \theta)$  to be an encoder that produces the content representation of an instance  $x$  with parameters  $\theta$ . Note that this encoder is shared between instances from both the clean and the weak set. Let  $f_c(enc(x); \gamma_c)$  and  $f_w(enc(\tilde{x}); \gamma_w)$  be the functions that map the content representation of the instances to their labels on the clean and weakly supervised data, respectively. Note that in contrast to the encoder with shared parameters  $\theta$ , the parameters  $\gamma_c$  and  $\gamma_w$  are different for the clean and weak sources respectively to capture their individual characteristics. The final objective for jointly optimizing the predictions from dual sources of supervision is given by:

$$\min_{\theta, \gamma_c, \gamma_w} \mathbb{E}_{(x, y) \in \mathcal{D}} \mathcal{L}(y, f_c(enc(x))) + \alpha \mathbb{E}_{(\tilde{x}, \tilde{y}) \in \tilde{\mathcal{D}}} \mathcal{L}(\tilde{y}, f_w(enc(\tilde{x}))) \quad (1)$$

where  $\mathcal{L}$  denotes the loss function to minimize the prediction error of the model.  $\alpha$  is a hyper-parameter that controls the relative importance of the loss functions computed over the data from clean and weak sources.

**3.1.1 Weak Label Correction.** Labeling functions are heuristic and can generate false labels (refer to Section 2.3.2 for an evaluation of labeling functions). An intuitive approach is to consider correcting these noisy labels before feeding them into the above model. Label correction methods have been previously studied for learning from noisy data sources [19, 45]. We now give a brief primer on existing work on label correction before discussing on how to integrate it into our framework.

*Prior Work on Label Correction.* Leveraging weak supervision for building effective supervised models has shown performance improvement in various tasks [19, 45]. Hendrycks *et al.* propose the idea of learning a *label corruption matrix* to estimate clean labels from the weak labels with the Gold Loss Correction approach (GLC) [19]. Given a set of instances  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$  with manually annotated (clean) labels  $y$  for  $L$  categories, and a weak labeled set  $\tilde{\mathcal{D}} = \{\tilde{x}_j, \tilde{y}_j\}_{j=1}^N$ , GLC aims to estimate a matrix  $C \in \mathcal{R}^{L \times L}$  to model the label corruption process. Formally, we first train a

<sup>3</sup>Hydra is a multi-headed creature in the Greek legend. Our proposed framework leverages multiple sources of supervision jointly and therefore is named Hydra.

**Table 4: Notation Table.**

Notation	Meaning
$\mathcal{D}$	set of instances with clean labels
$\tilde{\mathcal{D}}$	set of instances with weak labels
$\mathcal{D}'$	set of instances with weak label correction
$enc(\cdot)$	shared encoder for learning latent representations
$C$	label corruption matrix
$f_c$	function for predicting clean labels
$f_w$	function for predicting weak labels
$f'$	function for correcting weak labels
$w$	neural network model parameters
$v$	latent variable to select weak training samples

classifier  $f$  on the weakly labeled data  $\tilde{\mathcal{D}} = \{\tilde{x}_j, \tilde{y}_j\}_{j=1}^N$  as:

$$f(\tilde{x}) = \hat{p}(\tilde{y}|\tilde{x}, \theta)$$

Let  $\mathcal{X}_l$  be the subset of  $x$  with label  $y = l$ . Assuming the conditional independence of  $\tilde{y}$  and  $y$  given  $x$ , i.e.,  $p(\tilde{y}|y, x) = p(\tilde{y}|x)$ , we then estimate the corruption matrix  $\tilde{C}$  as follows,

$$C_{lr} = \frac{1}{|\mathcal{X}_l|} \sum_{x \in \mathcal{X}_l} \hat{p}(\tilde{y} = r|x) = \frac{1}{|\mathcal{X}_l|} \sum_{x \in \mathcal{X}_l} \hat{p}(\tilde{y} = r|y = l, x) \approx p(\tilde{y} = r|y = l) \quad (2)$$

With the new estimated  $C$ , they train a new classification model  $f'(x) = \tilde{p}(y|x, \theta)$  solving the following optimization problem:

$$\min_{\theta} \mathbb{E}_{(x, y) \in \mathcal{D}} \mathcal{L}(y, f'(x)) + \mathbb{E}_{(\tilde{x}, \tilde{y}) \in \tilde{\mathcal{D}}} \mathcal{L}(\tilde{y}, C^T f'(\tilde{x})) \quad (3)$$

where  $\mathcal{L}$  is a differentiable loss function to measure the prediction error, such as the cross-entropy loss.

**3.1.2 Integration with Hydra.** We use a similar idea and correct the weak labels for instances in  $\tilde{\mathcal{D}}$ . Using Equation 3, we learn a label correction function  $f'(\tilde{x})$  that rectifies the weak labels coming from the labeling functions for each instance  $\tilde{x} \in \tilde{\mathcal{D}}$ . We now obtain a label corrected weak supervision set  $\mathcal{D}' = \{\tilde{x}_j, f'(\tilde{x}_j)\}_{j=1}^N$ . Note that the label correction network reduces noise but the rectified labels could still be erroneous, and therefore considered as another source of weak supervision. In the new setting, we first feed the weakly labeled instances  $\tilde{\mathcal{D}}$  from the labeling function into the label correction network to obtain the rectified instances  $\mathcal{D}'$ . These are used as an input to Hydra. Formally, the overall objective function of our final model Hydra is given by:

$$\min_{\theta, \gamma_c, \gamma_w} \mathbb{E}_{(x, y) \in \mathcal{D}} \mathcal{L}(y, f_c(enc(x))) + \alpha \mathbb{E}_{(\tilde{x}, f'(\tilde{x})) \in \mathcal{D}'} \mathcal{L}(f'(\tilde{x}), f_w(enc(\tilde{x}))) \quad (4)$$

## 3.2 Self-paced Learning for Hydra

A simple training process is to consider all the weakly labeled samples jointly for learning. However, not all training samples are created equal. Some of the weak instances are noisier than others; whereas some are quite different in nature than the clean samples and therefore more difficult to learn from.

This is similar to curriculum learning [5] where a training schedule is used to first learn from easy samples followed by difficult ones.

The main challenge however is the distinction between easy and hard training samples. To alleviate this challenge, we can leverage the learned model to identify an easy set of samples given by a good fit in the model space similar to self-paced learning [27].

Consider  $v(\tilde{x}) \in \{0, 1\}$  to be a latent variable for each weak instance  $\tilde{x}$  that dictates whether to consider it for training. Correspondingly, our objective function is modified as follows.

$$\min_{\theta, \gamma_c, \gamma_w, v \in \{0, 1\}^N} \mathbb{E}_{(x, y) \in \mathcal{D}} \mathcal{L}(y, f_c(\text{enc}(x))) + \alpha \mathbb{E}_{(\tilde{x}, f'(\tilde{x})) \in \mathcal{D}'} [v(\tilde{x}) \cdot \mathcal{L}(f'(\tilde{x}), f_w(\text{enc}(\tilde{x}))) - \lambda \|v\|_1] \quad (5)$$

There are two distinct sets of parameters to learn corresponding to  $w = \{\theta, \gamma_c, \gamma_w\}$  for the neural network parameters and latent variables  $v$  for the training sample selection. To optimize the above equation, we employ alternate minimization. We first fix  $v$  and estimate the model parameters  $w$  using gradient descent.

Next we fix  $w$  and estimate  $v(\tilde{x})$  for all  $\tilde{x} \in \tilde{\mathcal{D}}$ . Partial derivative of Eqn. 5 with respect to  $v(\tilde{x})$  is given by  $\alpha \mathcal{L}(f'(\tilde{x}), f_w(\text{enc}(\tilde{x}))) - \lambda$ . The optimal solution for the equation is given by:

$$v(\tilde{x}) = \begin{cases} 1, & \text{if } \mathcal{L}(f'(\tilde{x}), f_w(\text{enc}(\tilde{x}))) < \frac{\lambda}{\alpha} \\ 0, & \text{otherwise} \end{cases}$$

Here  $\frac{\lambda}{\alpha}$  indicates whether an instance is easy to learn given by a small value of the corresponding loss function  $\mathcal{L}(\cdot)$ . A high loss indicates a poor fit of the sample in the model space and therefore ignored during training.  $\lambda$  as a hyper-parameter lets us control the injection of weak samples in the training set: a very low value admits few whereas a very high value admits all samples.

We initially train Hydra on only the clean data for a few epochs to trace the corresponding model space. Thereafter, we incorporate the weakly labeled samples gradually by increasing  $\lambda \in \{0.1, 0.2, \dots\}$  till all samples are included in the training set.

### 3.3 Training Hydra

We adopt mini-batch gradient descent with Adadelta [54] optimizer to learn the parameters. Adadelta is an adaptive method which divides the learning rate by an exponentially decaying average, and is less sensitive to the initial learning rate. For ease of understanding, all the notations we use are summarized in Table 4.

We first train the GLC model to obtain the label corrected weak supervision set  $\mathcal{D}'$ . To this end, we train a classifier  $f$  on weak supervision data  $\tilde{\mathcal{D}}$  and estimate the label corruption matrix  $C$ . Thereafter, we train a new classifier  $f'$  with the corruption matrix on the weakly supervised data, and obtain the data with corrected weak labels  $\mathcal{D}'$ .

Next we train Hydra for a few epochs on the clean data to have an initial estimate of  $w$ . Given  $w$  and an initial value of  $\lambda$ , we compute loss for all the weak instances and include those with loss less than  $\frac{\lambda}{\alpha}$  in the training set. This is followed by re-estimating  $w$ . We iterate over these steps and gradually increase  $\lambda$  until all the samples are accounted for or the model stops improving. For inference, the label of an instance  $x$  is predicted by  $y = f_c(\text{enc}(x))$ .

## 4 EXPERIMENTS

In this section, we present the experiments to evaluate the effectiveness of Hydra.

### 4.1 Experimental Settings

**4.1.1 Datasets.** We primarily perform experiments on the Avocado email collection. We perform experiments on three different tasks (intents): request information, schedule meeting and promise action. Section 2.4 discusses, in details, the data annotation process to obtain the clean labels and harnessing user interactions to obtain the weakly labeled instances. In addition to Avocado, we also perform an experiment to show the generalizability of our approach for transfer to another domain, namely the email collection for Enron (discussed in Section 4.6). Table 3 shows the dataset statistics. The datasets are balanced with equal number of positive and negative instances. Note that Avocado is the only public email collection with available user interaction logs available.

**4.1.2 Evaluation metric:** We pose our task as a binary classification problem for each of the intents. We use accuracy as the evaluation metric. We report results on the test set with the model parameters picked with the best validation accuracy on the dev set (Table 3 shows the data splits). All runs are repeated for 5 times and the average is reported. We compare different methods and techniques at different values of the *clean data ratio* defined as:

$$\text{clean ratio} = \frac{\# \text{clean labeled samples}}{\# \text{clean labeled samples} + \# \text{weak labeled samples}}$$

Specifically, we report performance with two settings with different clean data ratios:

- **All:** The setting where all available clean labels are used. According to our dataset construction process in Section 2.4 this corresponds to clean ratio = 10%.
- **Tiny:** The extreme setting where clean labels are down-sampled to account for only 1% of all data samples (clean and weak) to demonstrate the impact of weak supervision when extremely small amounts of clean labels are available.

We use following encoders for learning content representations:

- **AvgEmb:** AvgEmb learns the representation with the average word embedding of all the words in the message text. This is also the base model used in GLC [19].
- **BiLSTM [17]:** BiLSTM is a bi-directional recurrent neural network that learns the long-term dependencies in the text comprising of a sequence of tokens. During encoding, the input to the BiLSTM are the word embedding vectors of the input text.

We truncate or pad the text to have a uniform length of 128 tokens. For both AvgEmb and BiLSTM, we use the pre-trained 300-dimensional GloVe [41] embeddings to initialize the embedding matrix, and fine-tune it during training. For BiLSTM, we use one hidden layer and set the number of the hidden states as 300, and take the last hidden state as the final feature representation. We also employ a fully connected layer with 300 hidden states after the BiLSTM network to capture the interaction across feature dimensions. We use cross-entropy loss as  $\mathcal{L}$  in all settings. For self-paced learning, we set the number of epochs to 10 for each run with a specific  $\lambda$  varying from  $[0.1, 0.2, \dots]$ . We vary hyper-parameter  $\alpha \in \{0.1, 1, 10\}$  and choose the one that achieves the best performance on the validation set.

**4.1.3 Baselines:** We compare Hydra against the following baselines. The first set of baselines we consider is using the same base

**Table 5: Performance of the proposed approach compared to several baselines. Clean Ratio denotes the ratio of clean labels to all available labels (clean and weak) that is used to train the corresponding models. We show results for 10% (All) and 1% (Tiny) clean ratios. Hydra outperforms all the baselines in all settings.**

Clean Ratio (Setting)	Intent	Encoder (enc)	Clean	Weak	Clean + Weak	Pre-Weak	IWT	GLC	Hydra
10% (All)	RI	AvgEmb	0.649	0.523	0.616	0.613	0.661	0.693	<b>0.726</b>
		BiLSTM	0.688	0.524	0.684	0.717	0.711	0.717	<b>0.804</b>
	SM	AvgEmb	0.650	0.624	0.691	0.676	0.713	0.694	<b>0.731</b>
		BiLSTM	0.655	0.605	0.693	0.702	0.705	0.697	<b>0.714</b>
	PA	AvgEmb	0.641	0.628	0.633	0.637	0.625	0.647	<b>0.664</b>
		BiLSTM	0.608	0.547	0.611	0.631	0.616	0.635	<b>0.660</b>
1% (Tiny)	RI	AvgEmb	0.560	0.523	0.529	0.542	0.563	0.592	<b>0.664</b>
		BiLSTM	0.539	0.524	0.560	0.581	0.565	0.572	<b>0.622</b>
	SM	AvgEmb	0.565	0.624	0.618	0.633	0.628	0.620	<b>0.666</b>
		BiLSTM	0.538	0.605	0.626	0.608	0.625	0.617	<b>0.630</b>
	PA	AvgEmb	0.584	0.628	0.633	0.616	0.622	0.613	<b>0.647</b>
		BiLSTM	0.569	0.547	0.571	0.573	0.577	0.587	<b>0.626</b>

model as in Hydra: a three-layer neural network with word embeddings, an encoder (AvgEmb or BiLSTM) and a softmax layer for classification. We use this model in the following settings:

- Clean: Model trained on only the clean instances.
- Weak: Model trained on only the weak labels derived from user interactions. Weak labels are treated as regular clean labels.
- Clean+Weak: In this setting, we simply merge both the sets (essentially treating the weak labels to be as reliable as the clean ones) and use them together for training.
- Pre-Weak: We first pre-train the model on the weakly labeled instances. Then we take the trained model and fine-tune all the parameters in all the layers end-to-end on the clean instances.
- IWT: In Instance-Weighted Training (IWT), we assign sample weights to each of the instances during learning. For this, we modify Equation 1 as follows<sup>4</sup>.  

$$\min_{\theta, \gamma_c, \gamma_w} \mathbb{E}_{(x, y) \in \mathcal{D}} [u(x) \cdot \mathcal{L}(y, f_c(\text{enc}(x)))] + \alpha \mathbb{E}_{(\tilde{x}, \tilde{y}) \in \tilde{\mathcal{D}}} [v(\tilde{x}) \cdot \mathcal{L}(\tilde{y}, f_w(\text{enc}(\tilde{x})))]$$
with  $u > v$  forcing the model to focus more on the clean instances during learning. Note that the Clean+Weak baseline is a special case with  $u(x) = 1\forall x$ .

The next baseline is the Gold Loss Correction (GLC) [19] that estimates a label corruption matrix to model the correlation between weak and clean labels, which can be used to predict the unseen true labels. Finally, we report results from our full model Hydra.

## 4.2 Impact of Weak Supervision

Table 5 shows the performance comparing Hydra to other models in different settings. From therein we make the following observations based on varying clean ratio.

- Training only on the clean samples (even though they are much smaller in size) achieves better performance than training only on the weakly labeled ones on an aggregate across all the tasks and settings (demonstrated by Clean > Weak).
- Incorporating weakly labeled data even by simple aggregation with clean data like (Clean + Weak), pre-training (Pre-Weak) and

instance weighting (IWT) improves model performance on an aggregate over that of using only the clean or weak data.

- More sophisticated methods of integrating weakly labeled data with clean data gradually improves the performance on an aggregate across all the tasks (demonstrated by Hydra > GLC > IWT > Pre-Weak > Clean+Weak > Clean > Weak).
- Finally, irrespective of the clean ratio, Hydra achieves the best performance in all settings.

## 4.3 Impact of Clean Data Ratio

In real world scenarios, we often have a limited amount of annotated data, that would vary depending on the task and the domain, and a large amount of unlabeled data. In this experiment, we explore how the performance of Hydra changes with varying amount of clean data. To this end, we fix the number of weakly labeled instances and change the number of clean instances for each setting, thereby, varying the clean ratio between [1%, 3%, 5%, 7%, 9%, 10%]. Figure 3 shows the performance of different models using AvgEmb as the encoder. The graphs for BiLSTM encoder are similar and were omitted for space considerations. Since Clean + Weak, Pre-Weak and IWT have similar graphs, we only show results from the former. We make the following observations from Figure 3:

- With increasing clean ratio, the performance increases for all models (except Weak which uses a fixed amount of weakly labeled data). This is obvious as clean labels are more reliable than the weak ones.
- Hydra consistently outperforms all other methods accessing both clean and weak instances.
- Simple aggregation of clean and weak sources of supervision (e.g. Clean+Weak) without accounting for the source’s uncertainties (as modeled separately in Hydra) may not improve the prediction performance.
- The performance gap between using only the clean labels and Hydra using both clean and weak labels decreases with the increase in clean ratio. Note that Hydra is more effective when the clean ratio is small which is a more realistic setting with a small amount of labeled data.

<sup>4</sup>Results are reported with  $u(x) = 10\forall x$  and  $v(\tilde{x}) = 1\forall \tilde{x}$  with minimum error on the validation set.

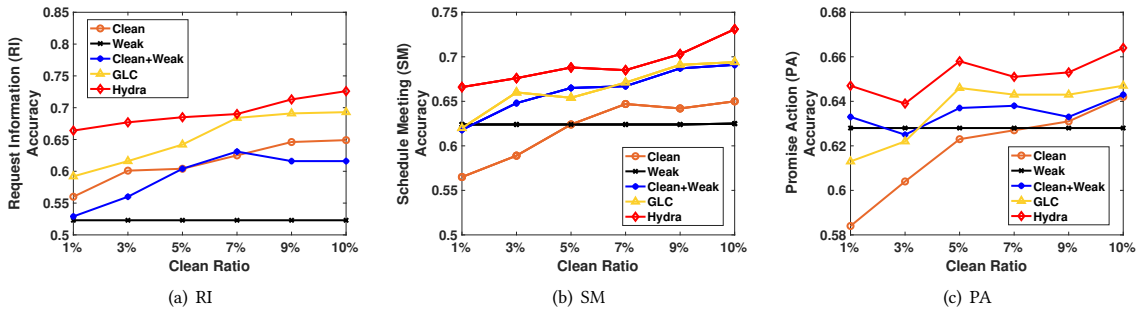


Figure 3: Classification results with varying clean ratio for different tasks ( $\text{enc}=\text{AvgEmb}$ ) (best viewed in color). We keep the number of weakly labeled instances fixed, and vary the amount of cleanly labeled ones for any given setting.

Table 6: Ablation analysis for Hydra. First row in each section shows Hydra with self-paced learning and GLC. Results are average across all tasks & encoders for a given clean ratio.

Clean Ratio (Setting)	Components	Accuracy
10% (All)	Hydra	<b>0.716</b>
	– self-paced learning	0.690
	– GLC	0.688
1% (Tiny)	Hydra	<b>0.643</b>
	– self-paced learning	0.631
	– GLC	0.632

#### 4.4 Ablation Analysis

**Self-paced Learning.** In order to understand the contribution of self-paced learning in the Hydra framework, we perform another experiment. In this we train Hydra leveraging all the clean and weak labels jointly without any curriculum. At each epoch, we sample batches with equal number of instances from clean and weakly labeled data and train Hydra end-to-end optimizing Equation 4. From Table 6, we observe self-paced learning to improve the performance of Hydra on aggregate across different tasks for different values of clean ratio. We observe the self-paced learning to perform much better at higher values of clean ratio that contributes a larger set of clean samples for training. This results from our training schedule (refer to Section 3.3) where we initially train Hydra for a few epochs on the clean data to trace the initial model space for  $w$ . **Gold Loss Correction (GLC).** In this, we remove the GLC component from Hydra and report the performance in Table 6. Similar to self-paced learning, we observe a similar performance loss at different values of the clean ratio. Note that both GLC and self-paced learning capture the noise in the weakly labeled instances to inform Hydra during training. However, the GLC component is learnt offline where Hydra uses only a fixed copy of the label corruption matrix. Whereas, self-paced learning updates all the parameters in Hydra during training.

#### 4.5 Additional Experiments

**Relative importance of clean and weak sources:** The hyperparameter  $\alpha$  controls the relative importance of the losses computed over the labels from the clean and weak sources. We observe that a

Table 7: Hydra ( $\text{enc} = \text{BiLSTM}$ ) on RI task with fixed amount of clean data (i.e. all the 1800 clean instances) and varying percentage of weak labels.

% weak labels	0%	20%	40%	60%	80%	100%
# weak labels	0	3,240	6,480	9,720	12,960	16,200
Accuracy	0.688	0.711	0.732	0.744	0.744	0.804

larger value of  $\alpha$  results in a better performance when the clean ratio is small, and vice versa when the clean ratio is large. This shows that Hydra relies more on the weak labels when the amount of clean labels is less. This dependency decreases with greater availability of reliable annotations. For example, when the clean ratio is 1% and 3% (see Figure 3), Hydra achieves the best performance with  $\alpha = 10$ ; whereas, with a clean ratio ranging from 5% to 10%,  $\alpha = 1$  leads to the best performance consistently across all the tasks.

**Impact of the amount of weak labels:** In Figure 3, we fixed the amount of weak labels and varied the amount of clean labels for different tasks depicting performance improvement with increase in the amount of reliable labels. In order to assess the impact of the amount of available weak labels on the performance of Hydra, we perform another experiment where we fix the task and the amount of clean labels, and vary the amount of weakly labeled instances. We perform this experiment with Hydra and BiLSTM encoder for the RI task. From Table 7, we observe that the performance of Hydra improves with increase in the amount of weak data, further demonstrating the importance of weak supervision.

#### 4.6 Domain Transfer

Now we want to test the transferability of Hydra. In other words, we want to test if the weak data from one domain can help in intent prediction in another domain. To this end, we train Hydra using the clean data in Enron and weak data in Avocado, and test the corresponding model on Enron. As baseline, we consider the base model used in Hydra trained using only the clean labeled data in Enron (see baseline description in Section 4.1 and data statistics in Table 3). All models have the same test set of 908 instances in Enron. Table 8 shows results for the task of scheduling meeting intent with different clean data ratio. We observe that Hydra trained on clean data from Enron and weak data from Avocado, and tested on Enron



**Table 8: Domain transfer for Hydra (enc=BiLSTM) on SM task. Av. → En. denotes Hydra trained on clean data in Enron and weak data in Avocado, and tested on Enron. Whereas En. → En. denotes the model trained on only the clean data in Enron, and tested on Enron.**

	En. → En.		Av. → En.	
# clean training labels	82	908	82	908
# weak training labels	0	0	8,176	8,176
Accuracy	0.714	0.717	0.752	0.821

shows better results than that trained only on the clean data from Enron. This shows that (i) Hydra transfers well across domains, and (ii) weak signals from one domain can be leveraged to improve the performance of models in another domain transferred via shared encoders. We also perform an experiment, where we train Hydra on Avocado using all the available clean and weak instances, and test on Enron. In this zero-shot transfer setting (with no labeled training data requirement for Enron), Hydra obtains an accuracy of 0.738. Comparing this setting to the above settings in Table 8, we observe that (i) Hydra performs better than En. → En. (using only clean data from Enron) (0.717) demonstrating transferability, and (ii) worse than Av. → En. (using clean data from Enron and weak data from Avocado) (0.821) demonstrating the benefit of adding some target domain clean data.

## 5 RELATED WORK

In this section, we briefly review the related work on email intent detection, weak supervision and learning from user interactions in other applications such as web search.

### 5.1 Email Intent Classification

Email understanding and intent classification has attracted increasing attention recently. Dabbish *et al.* conduct a survey on 124 participants to characterize email usage [13]. They highlight four distinct uses of email intents like project management, information exchange, scheduling and planning, and social communication. Detecting user intents, especially action-item intents [7], can help service providers to enhance user experience. Recent research focuses on predicting actionable email intent from email contents [31, 51], and identify related user actions such as reply [52], deferral [47], re-finding [33]. Wang *et al.* model the contextual information in email text to identify sentence-level user intents. Lin *et al.* build a reparameterized recurrent neural network to model cross-domain information and identify actionable email intents. In a more finer-grained level, Lampster *et al.* [28] study the problem of detecting emails that contain intent of requesting information, and propose to segment email contents into different functional zones. More recently, Azarbyonad *et al.* [4] utilize domain adaptation for commitment detection in emails. They demonstrate superior performance using autoencoders to capture both feature- and sample-level adaptation across domains. In contrast to all these models trained on manually annotated clean labels, we develop a framework Hydra that leverages weak supervision signals from user interactions for intent classification in addition to a small amount of clean labels.

### 5.2 Learning with Weak Supervision

Most machine learning models rely on the scale of labeled data to achieve good performance where the presence of label noise [37] or adversarial noise [44] can cause a dramatic performance drop. Therefore, learning with noisy labels has been of great interest to the research community for various tasks [15, 34, 36]. Some of the existing works attempt to rectify the weak labels by incorporating a loss correction mechanism [40, 48]. Sukhbaatar *et al.* [48] introduce a linear layer to adjust the loss and estimate label corruption with access to the true labels [48]. Patrini *et al.* [40] utilize the loss correction mechanism to estimate a label corruption matrix without making use of clean labels. Other works consider the scenario where a small set of clean labels are available [8, 19, 30, 45]. For example, Veit *et al.* use human-verified labels and train a label cleaning network in a multi-label classification setting. Recent works also consider the scenario where weak signals are available from multiple sources [42, 43, 49] to exploit the redundancy as well as the consistency in the labeling information. We build on top of previous work in the area of learning from weak supervision. More specifically, we focus on an application (email intent detection) where weak supervision can be obtained from user interaction signals. Second, we focus on a setup where a small amount of clean labeled data is available and propose methods to combine it with larger datasets with weak labels to improve the overall performance. In addition, our work does not make any strong assumptions about the structure of the noise or depend on the availability of multiple weak sources to model corroboration.

### 5.3 Learning from User Interactions in Web Search

Modern web search engines have heavily exploited user interaction logs to train and improve search systems [2, 23]. In fact, previous work [18] showed that models using click behaviors are more predictive of goal success than using document relevance. Motivated by the success of deep learning methods, several studies have focused on developing deep ranking models for web search. Since such models require access to large amounts of training data, they opted for using click data [32, 35] or the output of an unsupervised ranking model, BM25, as a weak labels [14].

Our work is similar to this line of work in that they both try to leverage user interaction data to improve a machine learning system. They are also different in many ways. First, we focus on intent classification in email text as opposed to ranking. Additionally, we focus on methods that combine clean-labeled data and weak-labeled data as opposed to just using implicit feedback data like clicks. Finally while clicks were shown to be more accurate than other types of user interaction signals, they suffer from several types of biases (e.g. snippet as, positing bias, etc.). Understanding the difference between such user interaction signals and clicks is an interesting direction for future research.

## 6 CONCLUSIONS AND FUTURE WORK

In this paper, we leverage weak supervision signals from user interactions to improve intent detection for emails. We develop an end-to-end robust neural network model Hydra to jointly learn from a small amount of clean labels and a large amount of weakly

labeled instances derived from user interactions. Extensive experiments on a real-world email dataset, Avocado, show Hydra to not only outperform state-of-the-art baselines but also its effectiveness in transferring the weak signals to another domain, namely Enron.

There are several directions for further investigation. First, we can extend our framework to multi-task learning where all of the above intent classification tasks can be learned jointly along with multiple sources of weak supervision. It is also worth exploring combining label correction and multi-source learning jointly instead of a two-stage approach. Second, understanding the nature of different sources of weak supervision is valuable for learning in different application domains. For example, in web search, user clicks can be a relatively accurate source of weak supervision, but may suffer from presentation bias; while for email data, user interactions are less accurate but may not suffer from the same biases.

## REFERENCES

- [1] Email statistics report. The Radicati Group, INC., 2015.
- [2] Eugene Agichtein, Eric Brill, and Susan Dumais. Improving web search ranking by incorporating user behavior information. In *SIGIR*, 2006.
- [3] Qingyao Ai, Susan T Dumais, Nick Craswell, and Dan Liebling. Characterizing email search using large-scale behavioral logs and surveys. In *WWW*, 2017.
- [4] Hosein Azarbyonad, Robert Sim, and Ryan W White. Domain adaptation for commitment detection in email. In *WSDM*, 2019.
- [5] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, 2009.
- [6] Paul N. Bennett and Jaime Carbonell. Detecting action-items in e-mail. In *SIGIR*, 2005.
- [7] Paul N Bennett and Jaime G Carbonell. Detecting action items in email. 2005.
- [8] Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In *SIGACT*, 2017.
- [9] Michael Chui, James Manyika, Jacques Bughin, Richard Dobbs, Charles Roxburgh, Hugo Sarrazin, Georey Sands, and Magdalena Westergren. The social economy: Unlocking value and productivity through social technologies. McKinsey Global Institute, 2012.
- [10] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.
- [11] William W. Cohen, Vitor R. Carvalho, and Tom M. Mitchell. Learning to classify email into speech acts. In *In Proceedings of Empirical Methods in Natural Language Processing*, 2004.
- [12] Laura A. Dabbish, Robert E. Kraut, Susan Fussell, and Sara Kiesler. Understanding email use: Predicting action on a message. In *CHI*. ACM, 2005.
- [13] Laura A Dabbish, Robert E Kraut, Susan Fussell, and Sara Kiesler. Understanding email use: predicting action on a message. In *CHI*, 2005.
- [14] Mostafa Dehghani, Hamed Zamani, Aliaksei Severyn, Jaap Kamps, and W Bruce Croft. Neural ranking models with weak supervision. In *SIGIR*, 2017.
- [15] Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2013.
- [16] Liang Ge, Jing Gao, Xiaoyi Li, and Aidong Zhang. Multi-source deep learning for information trustworthiness estimation. In *KDD*, 2013.
- [17] Alex Graves and Jürgen Schmidhuber. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural networks*, 18(5-6):602–610, 2005.
- [18] Ahmed Hassan, Rosie Jones, and Kristina Lisa Klinkner. Beyond dcg: user behavior as a predictor of a successful search. In *WSDM*, 2010.
- [19] Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. In *NeurIPS*, 2018.
- [20] Eric Horvitz. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '99, pages 159–166, New York, NY, USA, 1999. ACM.
- [21] Eric Horvitz, Andy Jacobs, and David Hovel. Attention-sensitive alerting. In *UAI*, 1999.
- [22] Meng Jiang, Peng Cui, Rui Liu, Qiang Yang, Fei Wang, Wenwu Zhu, and Shiqiang Yang. Social contextual recommendation. In *CKM*, 2012.
- [23] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, Filip Radlinski, and Geri Gay. Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search. *TOIS*, 25(2):7, 2007.
- [24] Anjali Kannan, Karol Kurach, Sujith Ravi, Tobias Kaufmann, Andrew Tomkins, Balint Miklos, Greg Corrado, László Lukács, Marina Ganea, Peter Young, et al. Smart reply: Automated response suggestion for email. *arXiv preprint arXiv:1606.04870*, 2016.
- [25] Bryan Klimt and Yiming Yang. The enron corpus: A new dataset for email classification research. In *ECML*, 2004.
- [26] Farshad Kooti, Luca Maria Aiello, Mihajlo Grbovic, Kristina Lerman, and Amin Mantrach. Evolution of conversations in the age of email overload. In *WWW*, 2015.
- [27] M Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *NeurIPS*, 2010.
- [28] Andrew Lampert, Robert Dale, and Cecile Paris. Detecting emails containing requests for action. In *HLT*, 2010.
- [29] Andrew Lampert, Robert Dale, and Cecile Paris. Detecting emails containing requests for action. In *NAACL*, 2010.
- [30] Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. Learning from noisy labels with distillation. In *ICCV*, 2017.
- [31] Chu-Cheng Lin, Dongyeop Kang, Michael Gamon, and Patrick Pantel. Actionable email intent modeling with reparametrized rnn. In *AAAI*, 2018.
- [32] Cheng Luo, Yukun Zheng, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. Training deep ranking model with weak relevance labels. In Zi Huang, Xiaokui Xiao, and Xin Cao, editors, *Databases Theory and Applications*, pages 205–216, Cham, 2017. Springer International Publishing.
- [33] Joel Mackenzie, Kshitiz Gupta, Fang Qiao, Ahmed Hassan Awadallah, and Milad Shokouhi. Exploring user behavior in email re-finding tasks. In *WWW*, 2019.
- [34] Yu Meng, Jiaming Shen, Chao Zhang, and Jiawei Han. Weakly-supervised hierarchical text classification. In *AAAI*, 2019.
- [35] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. Learning to match using local and distributed representations of text for web search. In *WWW*, 2017.
- [36] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *NeurIPS*, 2013.
- [37] David F Nettleton, Albert Orriols-Puig, and Albert Fornells. A study of the effect of different types of noise on the precision of supervised learning techniques. *Artificial intelligence review*, 33(4):275–306, 2010.
- [38] Douglas Oard, William Webber, David Kirsch, and Sergey Golitsynskiy. Avocado research email collection. *Philadelphia: Linguistic Data Consortium*, 2015.
- [39] Wanli Ouyang, Xiao Chu, and Xiaogang Wang. Multi-source deep learning for human pose estimation. In *CVPR*, 2014.
- [40] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, 2017.
- [41] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.
- [42] Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. Snorkel: Rapid training data creation with weak supervision. *VLDB*.
- [43] Alexander Ratner, Braden Hancock, Jared Dunmon, Frederic Sala, Shreyash Pandey, and Christopher Ré. Training complex models with multi-task weak supervision. *arXiv preprint arXiv:1810.02840*, 2018.
- [44] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*, 2014.
- [45] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. *arXiv preprint arXiv:1803.09050*, 2018.
- [46] Maya Sappelli, Gabriella Pasi, Suzan Verberne, Maaik de Boer, and Wessel Kraaij. Assessing e-mail intent and tasks in e-mail messages. *Information Sciences*, 358:1–17, 2016.
- [47] Bahareh Sarrafzadeh, Ahmed Hassan Awadallah, Christopher H Lin, Chia-Jung Lee, Milad Shokouhi, and Susan T Dumais. Characterizing and predicting email deferral behavior. In *WSDM*, 2019.
- [48] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2014.
- [49] Paroma Varma, Frederic Sala, Ann He, Alexander Ratner, and Christopher Ré. Learning dependency structures for weak supervision models. *ICML*, 2019.
- [50] Wei Wang, Saghar Hosseini, Ahmed Hassan Awadallah, Paul N. Bennett, and Chris Quirk. Context-aware intent identification in email conversations. In *SIGIR*, 2019.
- [51] Wei Wang, Saghar Hosseini, Ahmed Hassan Awadallah, Paul N Bennett, and Chris Quirk. Context-aware intent identification in email conversations. In *SIGIR*, 2019.
- [52] Liu Yang, Susan T Dumais, Paul N Bennett, and Ahmed Hassan Awadallah. Characterizing and predicting enterprise email reply behavior. In *SIGIR*, 2017.
- [53] Xiao Yang, Ahmed Hassan Awadallah, Madian Khabsa, Wei Wang, and Miaosen Wang. Characterizing and supporting question answering in human-to-human communication. In *SIGIR*, 2018.
- [54] Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- [55] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization, 2016.