# Prefer to Classify: Improving Text Classifiers via Auxiliary Preference Learning

**Jaehyung Kim**[1] Jinwoo Shin[1] Deogyeop Kang[2]

[1]Korea Advanced Institute of Science and Technology (KAIST)
[2]University of Minnesota (UMN)

# Importance of NLP Benchmarks

- Success of NLP systems has been driven by **large human-annotated benchmarks**
  - They guide the researchers in a right direction to develop methods
  - E.g., SQuAD (QA), GLUE (language understanding), and BIG-bench (large language models)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".
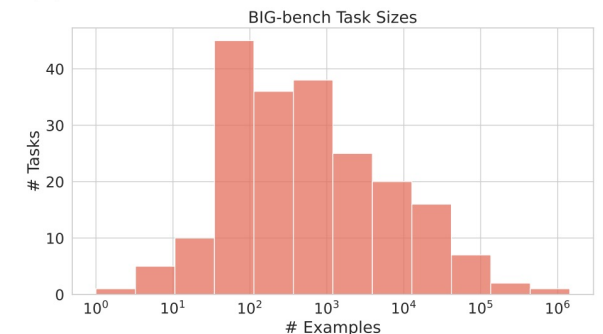
What causes precipitation to fall?
gravity

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
graupel

Where do water droplets collide with ice crystals to form precipitation?
within a cloud

| Corpus | \|Train\| | \|Test\| | Task |
|--------|-----------|----------|------|
| CoLA | 8.5k | **1k** | acceptability |
| SST-2 | 67k | 1.8k | sentiment |
| MRPC | 3.7k | 1.7k | paraphrase |
| STS-B | 7k | 1.4k | sentence similarity |
| QQP | 364k | **391k** | paraphrase |
| MNLI | 393k | **20k** | NLI |
| QNLI | 105k | 5.4k | QA/NLI |
| RTE | 2.5k | 3k | NLI |
| WNLI | 634 | **146** | coreference/NLI |

*Example of SQuAD (100k+)* [Rajpurkar et al. 2016]    *Summary of GLUE* [Wang et al. 2019]    *Diversity/scale of BIG-bench* [Srivastava et al. 2022]
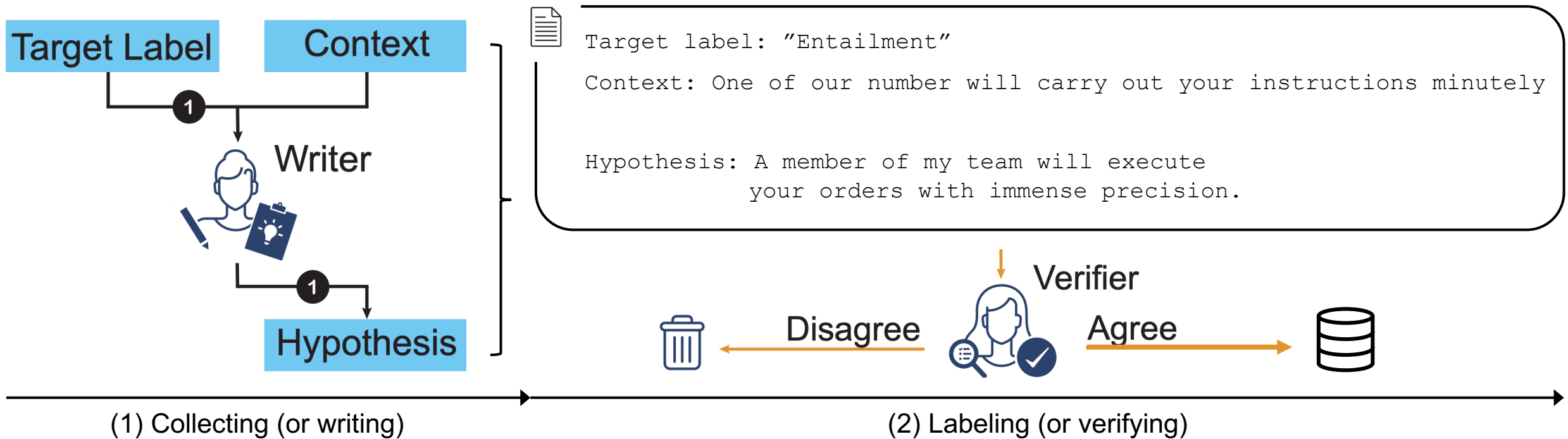
[Rajpurkar et al. 2016] SQuAD: 100,000+ Questions for Machine Comprehension of Text, EMNLP 2016
[Wang et al. 2019] GLUE: A Multi-task Benchmark and Analysis Platform for Natural Language Understanding, ICLR 2019
[Srivastava et al. 2022] Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models, arXiv:2206.04615
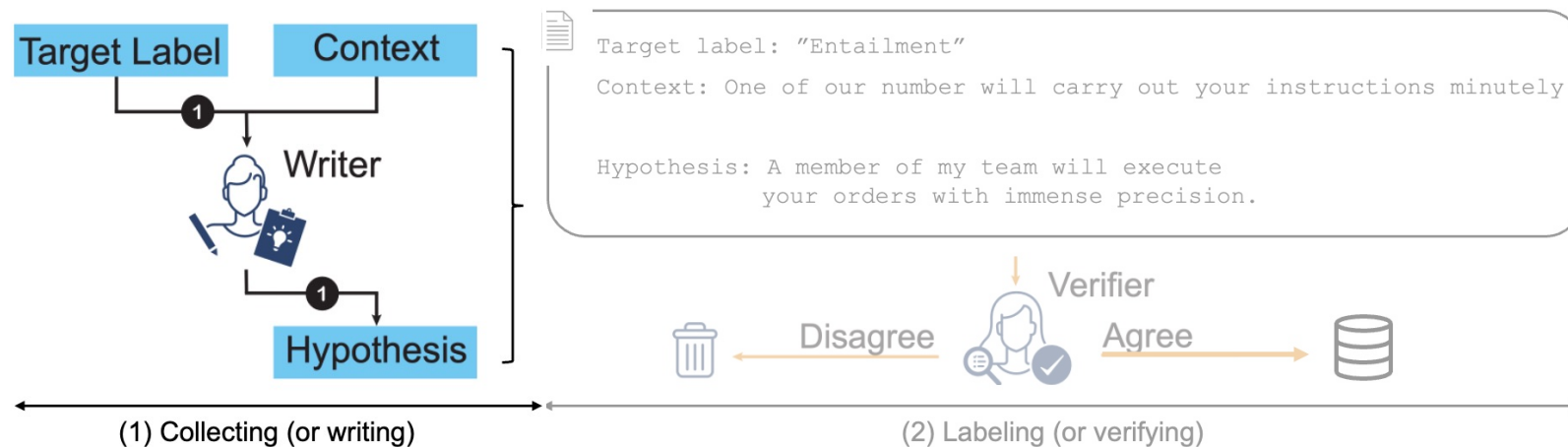
# Construction of NLP Benchmarks

- These benchmarks are usually constructed by following steps
  1. **Collecting (or writing)** the relevant input texts
  2. **Labeling** input texts **(or verifying)** by human annotators



Target label: "Entailment"

Context: One of our number will carry out your instructions minutely

Hypothesis: A member of my team will execute
your orders with immense precision.

(1) Collecting (or writing)

(2) Labeling (or verifying)

*Example of procedure for constructing benchmark for NLI task* [Nie et al. 2019]

[Nie et al. 2019] Adversarial NLI: A New Benchmark for Natural Language Understanding, ACL 2019

# Cost for Constructing NLP Benchmarks

- These benchmarks are usually constructed by following steps

  1. **Collecting (or writing)** the relevant input texts → **more costly and cumbersome**
     - E.g., distribution shift or spurious patterns of input make model suffer being generalized [Gururangan et al. 2018; Karamcheti et al. 2021]
     - Hence, much higher cost is often paid to the collection process to keep the quality [Kaushik et al. 2020]

  2. **Labeling** input texts **(or verifying)** by human annotators



Target Label    Context

Writer

Hypothesis

(1) Collecting (or writing)

Target label: "Entailment"

Context: One of our number will carry out your instructions minutely

Hypothesis: A member of my team will execute your orders with immense precision.

Verifier

Disagree    Agree

(2) Labeling (or verifying)

[Gururangan et al. 2018] Annotation Artifacts in Natural Language Inference Data, NAACL 2018
[Karamcheti et al. 2021] Mind Your Outliers! Investigating the Negative Impact of Outliers on Active Learning for Visual Question Answering, ACL 2021
[Kaushik et al. 2020] Learning the Difference that Makes a Difference with Counterfactually-augmented Data, ICLR 2020

# Complementary Way to Annotate Existing Benchmarks

- Hence, it is preferable to pay additional human cost to **auxiliary annotation**
  - E.g., improving label quality with more annotators [Nie et al. 2020]
  - or obtaining finer task information with new label space [Williams et al. 2020]

| Context | Hypothesis | Old Labels | New Labels majority and individual labels | Source | Type |
|---|---|---|---|---|---|
| With the sun rising, a person is gliding with a huge parachute attached to them. | The person is falling to safety with the parachute | Entailment E E E N N | Entailment $E^{(50)}$ $N^{(50)}$ | SNLI | Low agreements |
| A woman in a tan top and jeans is sitting on a bench wearing headphones. | A woman is listening to music. | Entailment E E N N E | Neutral $N^{(93)}$ $E^{(7)}$ | SNLI | Majority changed |
| A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match. | The men didn't appreciate the figure of the blonde woman sitting at the bar. | Contradiction C N N C C | Contradiction $C^{(56)}$ $N^{(44)}$ | MNLI | Low agreements |
| In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols. | He had spotted Adrin preparing to fire his pistols. | Neutral N E N N E | Entailment $E^{(94)}$ $N^{(5)}$ $C^{(1)}$ | MNLI | Majority changed |

*Analysis of existing NLI datasets with more annotations* [Nie et al. 2020]

| Dataset | Subset | Numerical | Basic | Reference | Tricky | Reasoning | Error |
|---|---|---|---|---|---|---|---|
| **A1** | All | 40.8 | 31.4 | 24.5 | 29.5 | 58.4 | 3.3 |
| | C | 18.6 | 8.2 | 7.8 | 13.7 | 11.9 | 0.7 |
| | N | 7.0 | 9.8 | 7.1 | 6.4 | 31.3 | 1.0 |
| | E | 15.2 | 13.4 | 9.6 | 9.4 | 15.2 | 1.6 |
| **A2** | All | 38.5 | 41.2 | 29.4 | 29.1 | 62.7 | 2.5 |
| | C | 15.6 | 11.8 | 10.2 | 13.6 | 15.5 | 0.3 |
| | N | 8.1 | 12.8 | 9.1 | 7.4 | 30.0 | 1.4 |
| | E | 14.8 | 16.6 | 10.1 | 8.1 | 17.2 | 0.8 |
| **A3** | All | 20.3 | 50.2 | 27.5 | 25.6 | 63.9 | 2.2 |
| | C | 8.7 | 17.2 | 8.6 | 12.7 | 14.9 | 0.3 |
| | N | 4.9 | 13.1 | 8.2 | 4.6 | 30.1 | 1.0 |
| | E | 6.7 | 19.9 | 10.7 | 8.3 | 18.9 | 0.8 |

*Analysis of ANLI with fine-grained annotation* [Williams et al. 2020]

[Nie et al. 2020] What Can We Learn from Collective Human Opinions on Natural Language Inference Data?, EMNLP 2020
[Williams et al. 2020] ANLIzing the Adversarial Natural Language Inference Dataset, arXiv:2010.12729

# Complementary Way to Annotate Existing Benchmarks

- Hence, it is preferable to pay additional human cost to **auxiliary annotation**
  - E.g., improving label quality with more annotators [Nie et al. 2020]
  - or obtaining finer task information with new label space [Williams et al. 2020]

| Context | Hypothesis | Old Labels | New Labels majority and individual labels | Source | Type |
|---|---|---|---|---|---|
| With the sun rising, a person is gliding with a huge parachute attached to them. | The person is falling to safety with the parachute | Entailment E E E N N | Entailment E$^{(50)}$ N$^{(50)}$ | SNLI | Low agreements |
| A woman in a tan top and jeans is sitting on a bench wearing headphones. | A woman is listening to music. | Entailment E E N N E | Neutral N$^{(93)}$ E$^{(7)}$ | SNLI | Majority changed |
| A group of guys went out for a drink after work, and sitting at the bar was a real a 6 foot blonde with a fabulous face and figure to match. | The men didn't appreciate the figure of the blonde woman sitting at the bar. | Contradiction C N N C C | Contradiction C$^{(56)}$ N$^{(44)}$ | MNLI | Low agreements |
| In the other sight he saw Adrin's hands cocking back a pair of dragon-hammered pistols. | He had spotted Adrin preparing to fire his pistols. | Neutral N E N N E | Entailment E$^{(94)}$ N$^{(5)}$ C$^{(1)}$ | MNLI | Majority changed |

| Dataset | Subset | Numerical | Basic | Reference | Tricky | Reasoning | Error |
|---|---|---|---|---|---|---|---|
|  | All | 40.8 | 31.4 | 24.5 | 29.5 | 58.4 | 3.3 |
| **A1** | C | 18.6 | 8.2 | 7.8 | 13.7 | 11.9 | 0.7 |
|  | N | 7.0 | 9.8 | 7.1 | 6.4 | 31.3 | 1.0 |
|  | E | 15.2 | 13.4 | 9.6 | 9.4 | 15.2 | 1.6 |
|  | All | 38.5 | 41.2 | 29.4 | 29.1 | 62.7 | 2.5 |
| **A2** | C | 15.6 | 11.8 | 10.2 | 13.6 | 15.5 | 0.3 |
|  | N | 8.1 | 12.8 | 9.1 | 7.4 | 30.0 | 1.4 |
|  | E | 14.8 | 16.6 | 10.1 | 8.1 | 17.2 | 0.8 |
|  | All | 20.3 | 50.2 | 27.5 | 25.6 | 63.9 | 2.2 |
| **A3** | C | 8.7 | 17.2 | 8.6 | 12.7 | 14.9 | 0.3 |
|  | N | 4.9 | 13.1 | 8.2 | 4.6 | 30.1 | 1.0 |
|  | E | 6.7 | 19.9 | 10.7 | 8.3 | 18.9 | 0.8 |

🧐 Can we find a **new alternative way to better exploit** existing <u>benchmarks</u> (input texts and task labels) via **auxiliary annotation?**
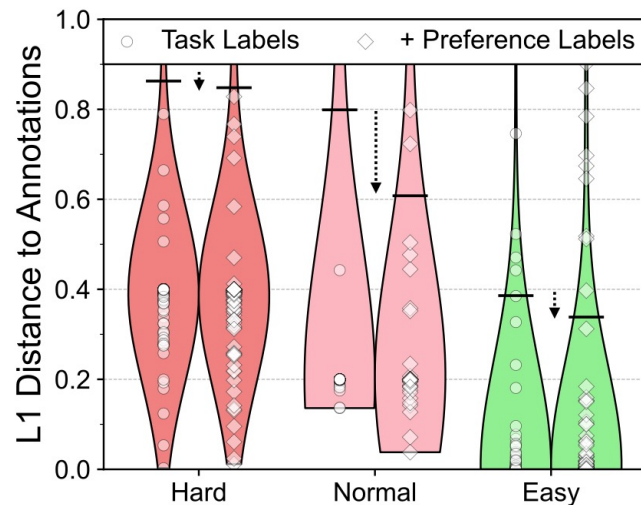Especially, for text classification
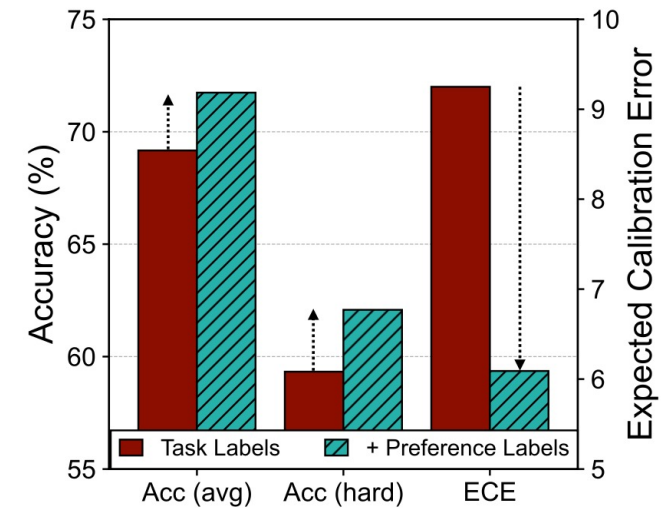
# Task-specific Preference as Auxiliary Annotation

- Idea: using **task-specific preference** between input texts as auxiliary annotation
  - To improve the text classification system upon existing task annotations
  - Auxiliary preference learning provides additional informative training signal to model
    - By relatively ordering a pair of two texts and better calibrating them w.r.t task through "pair-wise" comparison
    
    v.s. "instance-wise" task annotation



(a) Pair-wise preference signals

(b) Alignment to human annotations

(c) Improved text classification

*Concept of auxiliary preference learning and its empirical advantages*

# Prefer to Classify (P2C)

- Specifically, we propose following components for auxiliary preference learning
  - **Three different types of preference labels** in practical scenario
    - Using large language model (generative), data annotation records (extractive), or crowd workers (subjective)
  - **Novel multi-task learning** framework with task and preference labels: prefer-to-classify (P2C)



Multi-task learning with task and preference labels

Preference label with three different ways

*Visual illustration of the proposed auxiliary preference learning for improving text classifier*

4

# Different Types of Preference Labels

- **3 different types of preference** labels to apply auxiliary preference learning via P2C
  - **Generative preference** from large language models, e.g., GPT-3 [Brown et al. 2020]
    - <u>Good quality</u> from strong zero/few-shot generalization capability of LM, <u>low cost</u>, and **easy to access**



*Prompt design to collect generative preference labels from GPT-3* [Brown et al. 2020]

[Brown et al. 2020] Language Models are Few-shot Learners, NeurIPS 2020

# Different Types of Preference Labels

- **3 different types of preference** labels to apply auxiliary preference learning via P2C
  - **Extractive preference** from data annotation records
    - If one sample has higher voting than the other sample as specific label, then it is assumed to be more preferred
    - **Zero cost** with <u>good quality</u> by better exploiting information within task annotations, but often <u>hard to access</u>
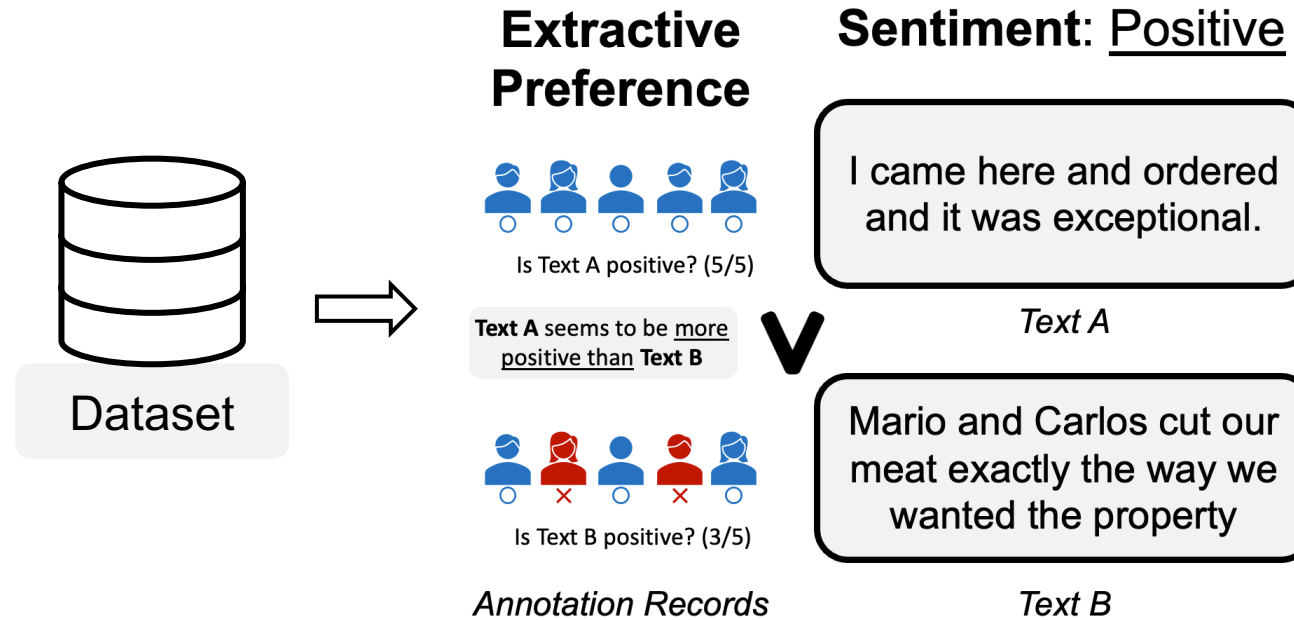


*Illustration of extractive preference from existing data annotation records*

# Different Types of Preference Labels

- **3 different types of preference** labels to apply auxiliary preference learning via P2C
  - **Subjective preference** from crowd workers
    - Obtained by directly asking the to humans, e.g., "*which sentence is more positive?*"
    - **Most accurate**, but it requires high cost and hence hard to access



*Used interface to collect subjective preference labels from crowd workers via AMT*

# Different Types of Preference Labels

- **3 different types of preference** labels to apply auxiliary preference learning via P2C
  - **Generative** / **Extractive** / **Subjective** preference labels
    - <u>Accuracy</u>: subjective > extractive ~= generative
    - <u>Cost</u>: extractive > generative >> subjective (e.g., 1.6$ for 10 subjective labels, while 8.0$ for 5,000 generative labels)
    - <u>Accessibility</u>: generative > extractive > subjective



**A**: We enjoyed our first and last meal in Toronto at Bombay Palace, and I can't think of a better way to book our journey. | **B**: So glad I finally tried this place because if confirmed my suspicions about that critic who rated it a 10.

Sentiment: <u>Positive</u>, Generative Preference: **A** ≻ **B**, Extractive Preference: **B** ≻ **A**, Subjective Preference: **No preference**

**A**: The buffalo chicken was not good, but very costly. | **B**: There was so much stuff from all over that I had to leave to find an ATM for more cash to pay for it all.

Sentiment: <u>Negative</u>, Generative Preference: **A** ≻ **B**, Extractive Preference: **B** ≻ **A**, Subjective Preference: **B** ≻ **A**

**A**: The hotel offered complimentary breakfast. | **B**: My friends had a full acrylic and the other had a fill. It looked so good.
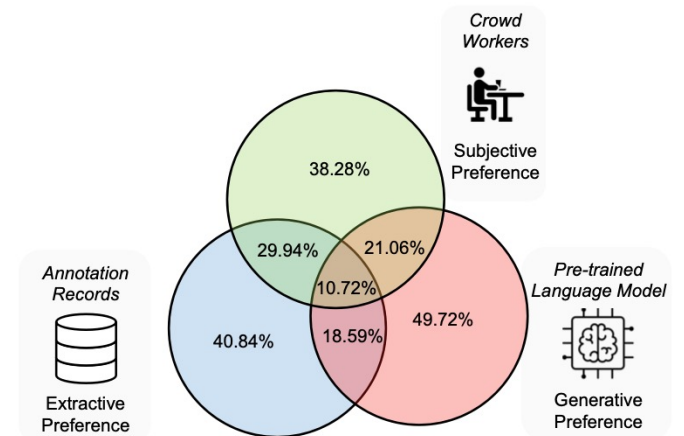
Sentiment: <u>Positive</u>, Generative Preference: **A** ≻ **B**, Extractive Preference: **A** ≻ **B**, Subjective Preference: **A** ≻ **B**

*Examples of the collected preference labels on same pair of sentences from DynaSent-R2* [Potts et al. 2021]

*Overlap between preferences*

[Potts et al. 2021] Dynasent: A Dynamic Benchmark for Sentiment Analysis, ACL 2021

# Prefer to Classify (P2C): Multi-task Learning

- Classifier is trained using task labels and preference labels jointly with
  - **Diverse multi-preference heads** for better preference modeling
    - For preference predictor, we add preference prediction head $W_{\text{pref}}$ on classifier $g_\phi(\mathbf{x})$ (e.g., BERT)

$$P_\psi[\mathbf{x}^1 \succ \mathbf{x}^0; y_{\text{task}}] = \frac{\exp\left(h_\psi(\mathbf{x}^1, y_{\text{task}})\right)}{\sum_{i \in \{0,1\}} \exp\left(h_\psi(\mathbf{x}^i, y_{\text{task}})\right)} \qquad h_\psi(\mathbf{x}, y_{\text{task}}) = W_{\text{pref}} \circ [g_\phi(\mathbf{x}); y_{\text{task}}]$$

$$\mathcal{L}_{\text{pref}} = - \mathop{\mathbb{E}}_{\substack{(\mathbf{x}^0, \mathbf{x}^1, y_{\text{task}}, y_{\text{pref}}) \\ \sim \mathcal{D}}} \left[ y_{\text{pref}} \log P_\psi[\mathbf{x}^1 \succ \mathbf{x}^0; y_{\text{task}}] + (1 - y_{\text{pref}}) \log P_\psi[\mathbf{x}^0 \succ \mathbf{x}^1; y_{\text{task}}] \right]$$

- Then, we introduce multiple preference heads $\{W_{\text{pref}}^{(t)}\}_{t=1}^T$ and maximize KL divergence between their prediction

$$\mathcal{L}_{\text{div}} = \frac{-1}{T-1} \sum_{j=1, j \neq i}^T D_{\text{KL}}\left(P_{\psi^{(i)}}(\mathbf{x}^1, \mathbf{x}^0; y_{\text{task}}) \| P_{\psi^{(j)}}(\mathbf{x}^1, \mathbf{x}^0; y_{\text{task}})\right)$$

- Overall multi-task learning objective

$$\mathcal{L}_{\text{multi}} = \mathcal{L}_{\text{task}} + \mathcal{L}_{\text{pref}}^{\text{all}} + \lambda_{\text{div}} \mathcal{L}_{\text{div}} \qquad \mathcal{L}_{\text{pref}}^{\text{all}} = \sum_{t=1}^T \mathcal{L}_{\text{pref}}^{\psi^{(t)}}$$

# Prefer to Classify (P2C): Multi-task Learning

- Classifier is trained using task labels and preference labels jointly with
  - **Diverse multi-preference heads** for better preference modeling
  - **Consistency regularization** between task and preference learning
    - To explicitly impose the intuition: "*preferred instance should have a higher confidence*"
    - To this end, applying following consistency loss

$$\mathcal{L}_{\text{cons}} = y_{\text{pref}} \max\{0, p_y(\mathbf{x}^1) - p_y(\mathbf{x}^0)\} + (1 - y_{\text{pref}}) \max\{0, p_y(\mathbf{x}^0) - p_y(\mathbf{x}^1)\}$$

    - Overall, our training loss is as follow:

$$\mathcal{L}_{\text{train}} = \mathcal{L}_{\text{multi}} + \lambda_{\text{cons}}\mathcal{L}_{\text{cons}}$$

# Prefer to Classify (P2C): Multi-task Learning

- Classifier is trained using task labels and preference labels jointly with
  - **Diverse multi-preference heads** for better preference modeling
  - **Consistency regularization** between task and preference learning
  - **Selecting informative pairs** of input texts
    - *Disagreement-based* sampling: selecting pairs with high variance across multiple preference predictors $\{h_{\psi^{(i)}}\}_{t=1}^{T}$
    - *Inconsistency-based* sampling: selecting pairs with high consistency loss $\mathcal{L}_{\text{cons}}$

# Experiments

- Text classification with **generative preference**
  - P2C is consistently effective in improving the performance (accuracy and calibration)
    - bAcc: balanced accuracy as datasets have imbalanced distribution, wAcc: worst-group accuracy (minority)
  - P2C also outperforms GPT-3 baselines → Not just distilling "instance-wise" knowledge of GPT-3

| Method | CoLA Mcc($\uparrow$) | CoLA ECE($\downarrow$) | SMS Spam bAcc($\uparrow$) / wAcc($\uparrow$) | SMS Spam ECE($\downarrow$) | Hate Speech bAcc($\uparrow$) / wAcc($\uparrow$) | Hate Speech ECE($\downarrow$) | Emotion bAcc($\uparrow$) / wAcc($\uparrow$) | Emotion ECE($\downarrow$) |
|---|---|---|---|---|---|---|---|---|
| Vanilla | $63.7_{\pm1.0}$ | $\underline{3.6}_{\pm1.6}$ | $96.9_{\pm0.3}$ / $\underline{95.1}_{\pm1.5}$ | $1.3_{\pm0.3}$ | $81.1_{\pm1.8}$ / $69.9_{\pm4.6}$ | $5.1_{\pm1.0}$ | $88.6_{\pm2.3}$ / $76.1_{\pm7.8}$ | $4.0_{\pm1.1}$ |
| Label Smoothing | $63.9_{\pm0.3}$ | $4.6_{\pm1.2}$ | $96.9_{\pm0.8}$ / $94.0_{\pm1.5}$ | $1.1_{\pm0.3}$ | $81.5_{\pm0.9}$ / $\underline{71.3}_{\pm3.2}$ | $6.6_{\pm1.0}$ | $\underline{89.8}_{\pm0.8}$ / $\underline{76.9}_{\pm6.6}$ | $4.0_{\pm0.9}$ |
| Max Entropy | $64.1_{\pm0.3}$ | $4.5_{\pm0.4}$ | $\underline{96.9}_{\pm1.1}$ / $94.7_{\pm1.6}$ | $1.2_{\pm0.3}$ | $\underline{81.6}_{\pm1.8}$ / $70.5_{\pm4.2}$ | $\underline{4.3}_{\pm0.7}$ | $89.1_{\pm1.1}$ / $73.1_{\pm2.5}$ | $\mathbf{3.6}_{\pm0.9}$ |
| CS-KD | $\underline{64.5}_{\pm1.4}$ | $4.1_{\pm1.1}$ | $96.8_{\pm0.9}$ / $94.0_{\pm2.4}$ | $1.1_{\pm0.2}$ | $81.4_{\pm2.6}$ / $69.6_{\pm5.1}$ | $5.3_{\pm1.8}$ | $89.4_{\pm1.6}$ / $74.0_{\pm6.8}$ | $4.1_{\pm0.2}$ |
| GPT-3 (0-shot) | $60.4$ | - | $90.3$ / $84.3$ | - | $68.7$ / $41.6$ | - | $50.2$ / $23.3$ | - |
| GPT-3 (5-shot) | $58.5_{\pm0.4}$ | - | $92.2_{\pm0.5}$ / $88.5_{\pm0.7}$ | - | $78.5_{\pm2.0}$ / $70.3_{\pm3.6}$ | - | $46.6_{\pm0.6}$ / $30.3_{\pm2.6}$ | - |
| GPT-3 (20-shot) | $58.3_{\pm1.4}$ | - | $95.8_{\pm0.4}$ / $94.4_{\pm0.7}$ | - | $77.8_{\pm0.5}$ / $69.0_{\pm1.5}$ | - | $47.5_{\pm1.0}$ / $30.8_{\pm4.5}$ | - |
| P2C (Ours) | $\mathbf{65.4}_{\pm1.0}$ | $\mathbf{2.8}_{\pm1.1}$ | $\mathbf{97.4}_{\pm0.4}$ / $\mathbf{95.2}_{\pm1.0}$ | $\mathbf{1.1}_{\pm0.3}$ | $\mathbf{82.4}_{\pm1.3}$ / $\mathbf{73.6}_{\pm4.5}$ | $\mathbf{4.0}_{\pm0.3}$ | $\mathbf{90.7}_{\pm0.7}$ / $\mathbf{81.7}_{\pm4.7}$ | $\mathbf{3.6}_{\pm0.8}$ |

**11.55%** relative test error reduction compared to *Vanilla*

*Test accuracy of fine-tuned RoBERTa-base classifiers*

# Experiments

- Text classification with **extractive preference (Free!)**
  - P2C even outperforms the strong baselines for learning with annotation records

| Method | Offensive | Polite-Wiki | Polite-SE | MNLI | DynaSent-R1 | DynaSent-R2 |
|---|---|---|---|---|---|---|
| Vanilla | $75.88_{\pm 0.72}$ | $89.35_{\pm 1.53}$ | $70.00_{\pm 1.49}$ | $81.92_{\pm 0.70}$ | $80.43_{\pm 0.30}$ | $71.23_{\pm 1.05}$ |
| Soft-labeling | $76.08_{\pm 1.44}$ | $89.57_{\pm 1.76}$ | $70.35_{\pm 1.68}$ | $\underline{82.67}_{\pm 0.50}$ | $81.10_{\pm 0.33}$ | $\underline{72.15}_{\pm 1.59}$ |
| Margin Loss | $\underline{76.67}_{\pm 1.18}$ | $88.51_{\pm 0.93}$ | $\underline{70.51}_{\pm 1.16}$ | $81.41_{\pm 0.63}$ | $80.42_{\pm 0.23}$ | $69.27_{\pm 0.98}$ |
| Filtering | $76.13_{\pm 1.18}$ | $89.50_{\pm 0.87}$ | $68.28_{\pm 2.43}$ | $82.13_{\pm 0.67}$ | $80.38_{\pm 0.34}$ | $69.86_{\pm 0.78}$ |
| Weighting | $76.17_{\pm 1.18}$ | $89.65_{\pm 1.46}$ | $68.38_{\pm 1.67}$ | $82.48_{\pm 0.49}$ | $80.21_{\pm 0.41}$ | $71.81_{\pm 1.12}$ |
| Multi-annotator | $76.50_{\pm 1.98}$ | $\underline{89.88}_{\pm 1.82}$ | $69.39_{\pm 2.84}$ | $82.61_{\pm 0.70}$ | $\underline{81.14}_{\pm 0.55}$ | $71.97_{\pm 1.25}$ |
| CS-KD | $75.75_{\pm 0.66}$ | $89.65_{\pm 1.84}$ | $70.10_{\pm 1.29}$ | $82.32_{\pm 0.23}$ | $80.63_{\pm 0.27}$ | $71.81_{\pm 0.67}$ |
| P2C (Ours) | $\mathbf{77.81}_{\pm 0.21}$ | $\mathbf{91.06}_{\pm 0.64}$ | $\mathbf{71.21}_{\pm 0.93}$ | $\mathbf{83.15}_{\pm 0.29}$ | $\mathbf{81.50}_{\pm 0.39}$ | $\mathbf{73.06}_{\pm 0.31}$ |

**7.59% / 4.27%** relative test error reduction compared to *Vanilla / Best,* respectively

*Test accuracy of fine-tuned RoBERTa-base classifiers*

# Experiments

- Comparison between different annotation methods
  - **Setup**: Given the existing datasets, adding the same number of annotations but different ways
  - **Results**
    - Overall, preference labels are effective for hard samples (i.e., high disagreement) along with strong calibration effects
    - Subjective preference labels are the most effective for improving accuracy and calibration

| Method | $N_{task}$ | $N_{pref}$ | $Acc_{avg}(\uparrow)$ | $Acc_{hard}$ / $Acc_{easy}(\uparrow)$ | $ECE(\downarrow)$ | $d_{hard}$ / $d_{easy}(\downarrow)$ |
|---|---|---|---|---|---|---|
| Vanilla | 7.5k | - | $69.03_{\pm1.29}$ | $59.33_{\pm2.57}$ / $80.00_{\pm1.22}$ | $9.25_{\pm1.39}$ | $0.856_{\pm0.01}$ / $0.405_{\pm0.03}$ |
| Task Labels | 12.5k | - | $71.17_{\pm1.35}$ | $57.86_{\pm2.31}$ / $\mathbf{84.21}_{\pm1.05}$ | $9.19_{\pm1.36}$ | $0.878_{\pm0.04}$ / $\mathbf{0.327}_{\pm0.02}$ |
| Generative Preference | 7.5k | 5k | $\underline{71.46}_{\pm1.16}$ | $\underline{61.77}_{\pm0.94}$ / $82.28_{\pm1.01}$ | $\underline{6.64}_{\pm0.79}$ | $0.850_{\pm0.02}$ / $0.361_{\pm0.02}$ |
| Extractive Preference | 7.5k | 5k | $71.36_{\pm1.19}$ | $61.16_{\pm1.91}$ / $\underline{83.11}_{\pm1.78}$ | $6.75_{\pm0.78}$ | $\underline{0.847}_{\pm0.03}$ / $\underline{0.351}_{\pm0.03}$ |
| Subjective Preference | 7.5k | 5k | $\mathbf{71.74}_{\pm1.04}$ | $\mathbf{62.08}_{\pm0.94}$ / $83.01_{\pm1.27}$ | $\mathbf{6.09}_{\pm0.31}$ | $\mathbf{0.828}_{\pm0.02}$ / $0.356_{\pm0.02}$ |

*Test accuracy of fine-tuned RoBERTa-base classifiers on DynaSent-R2*

# Summary

- We introduce **preference label** as new auxiliary annotation to improve benchmark
  - It provides additional informative training signal to model via "*pair-wise*" comparison
  - We propose an effective multi-task learning framework, coined *prefer-to-classify (P2C)*
  - We provide *three different ways* to obtain preference labels (generative/extractive/subjective)

- P2C shows <span style="color:blue">consistent improvements</span> on various NLP benchmarks
  - Improved test accuracy with better calibration

- P2C suggests <span style="color:blue">new way to evolve benchmark</span> along with recent advance of LM

For more details and results,
please see our paper and code

Thank you for attention 😄

arXiv: 2306.04925

Github