

# Paper Reading Report

Xiangqing Shen

xiangqing.shen@njust.edu.cn

Text Mining Lab (NUSTM)  
Nanjing University of Science and Technology

October 7, 2021

# Contents

- ❶ NEUROLOGIC DECODING
- ❷ TEMPORAL REASONING
- ❸ (COMET-)ATOMIC<sub>20</sub><sup>20</sup>
- ❹ Hashtags, Emotions, and Comments

# Contents

- ❶ NEUROLOGIC DECODING
- ❷ TEMPORAL REASONING
- ❸ (COMET-)ATOMIC<sub>20</sub><sup>20</sup>
- ❹ Hashtags, Emotions, and Comments

# NEUROLOGIC DECODING

## NEUROLOGIC DECODING: (Un)supervised Neural Text Generation with Predicate Logic Constraints

Ximing Lu<sup>†‡</sup> Peter West<sup>†‡</sup> Rowan Zellers<sup>†‡</sup>

Ronan Le Bras<sup>‡</sup> Chandra Bhagavatula<sup>‡</sup> Yejin Choi<sup>†‡</sup>

<sup>†</sup>Paul G. Allen School of Computer Science & Engineering, University of Washington

<sup>‡</sup>Allen Institute for Artificial Intelligence

{lux32, pawest, rowanz, yejin}@cs.washington.edu

{ronanlb, chandrab}@allenai.org

Figure 1: Paper information. [1]

[1] Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, Yejin Choi. NeuroLogic Decoding: (Un)supervised Neural Text Generation with Predicate Logic Constraints. NAACL2021.

# Task Definition

COMMONGEN (Lin et al., 2019)	
input	Concept-Set <span>food   table   sit   front</span> Constraints $(\text{food} \vee \text{foods}) \wedge (\text{table} \vee \text{tables}) \wedge$ $(\text{sit} \vee \text{sits} \vee \text{sat} \vee \text{sitting}) \wedge (\text{front} \vee \text{fronts})$
	Scenario The man <b>sat</b> with his <b>food</b> at the <b>front</b> of the <b>table</b> The <b>food</b> is in <b>front</b> of you <b>sit</b> at the <b>table</b> . a <b>table</b> of <b>food</b> <b>sits</b> in <b>front</b> of three people
Evaluate Gender Bias in MT (Stanovsky et al., 2019)	
input	Source <span>✓Bäckerin ✗Bäcker</span> The physician told the <b>baker</b> that she had cancer.
	Constraints <span><math>(\text{Ärztin} \vee \text{Arzt}) \wedge (\text{Bäckerin} \wedge \neg \text{Bäcker})</math></span> Target Der Arzt sagte dem <b>Bäckerin</b> , dass er Krebs habe.
Recipe Generation (Kiddon et al., 2016)	
input	Dish name <span>garlic butter steak</span> Ingredients 2 tsp <b>butter</b> , 1 <b>beef steak</b> , 1/4 tsp <b>soy sauce</b> , 1 tsp <b>parsley</b> , 1/8 tsp <b>salt</b> , 1/2 tsp <b>garlic</b> Constraints $\text{butter} \wedge (\text{beef} \vee \text{steak} \vee \text{meat}) \wedge \text{soy sauce} \wedge$ $(\text{parsley} \vee \text{herb}) \wedge \text{salt} \wedge (\text{garlic} \vee \text{vegetable}) \wedge$ $(\neg \text{pork} \wedge \neg \text{bean} \wedge \neg \dots) \leftarrow \text{any extra ingredients}$
	Recipe Mix 1 tablespoon <b>butter</b> , <b>parsley</b> , <b>garlic</b> and <b>soy sauce</b> . Sprinkle <b>steak</b> with <b>salt</b> . In a large skillet, heat remaining <b>butter</b> over medium heat. Add <b>steak</b> ; cook until <b>meat</b> reaches desired doneness, 4-7 minutes per side. Serve with <b>garlic butter</b> .

Figure 2: An example for conditional text generation.

# Limitations of Previous Work

- **Finetuning LMs on a dataset of task-specific examples.** However, PLMs struggle at learning to follow these constraints.
- **Mismatch caused by a fundamental under-specification of finetuning.** Improvements come from **constrained generation or learning the language style?** When increasing the finetuning data fed to GPT2 by **an order of magnitude**, constraint-satisfaction with standard beam search shows **only modest improvement**.

# Contributions

- Proposing **NEUROLOGIC DECODING**, which effectively enforces the satisfaction of given lexical constraints by controlling the decoding stage of sequence generation.
- **Converting the hard logic constraints into a soft penalty term in the decoding objective**, and use a beam-based search to find approximately-optimal solutions.
- Empirical results demonstrate that NEUROLOGIC DECODING ensures the satisfaction of given constraints while maintaining high generation quality, in turn leading to **new SOTA results in both the supervised and zero-shot setting**.

# Prerequisite

## Predicate $D(\mathbf{a}, \mathbf{y})$

Let us define a predicate  $D(\mathbf{a}, \mathbf{y})$  to be a boolean function indicating the occurrence of key phrase  $\mathbf{a}$  in a sequence  $\mathbf{y}$ , where  $\mathbf{a}$  can be either unigram or multi-gram.  $D(\mathbf{a}, \mathbf{y})$  will be true iff  $\mathbf{a}$  occurs in  $\mathbf{y}$ .

$$D(\mathbf{a}, \mathbf{y}) \equiv \exists i, \mathbf{y}_{i:i+|\mathbf{a}|} = \mathbf{a}$$

NEUROLOGIC accepts lexical constraints in Conjunctive Normal Form:

$$\underbrace{(D_1 \vee D_2 \cdots \vee D_i)}_{C_1} \wedge \cdots \wedge \underbrace{(D_k \vee D_{k+1} \cdots \vee D_n)}_{C_m}$$

## Notation

- Each individual constraint  $D_i \rightarrow$  a *literal*.
- The disjunction of literals  $\rightarrow$  a *clause*, denoted as  $C_j$ .



# Prerequisite

## Objective

The method seeks optimal sequences in which all clauses are satisfied:

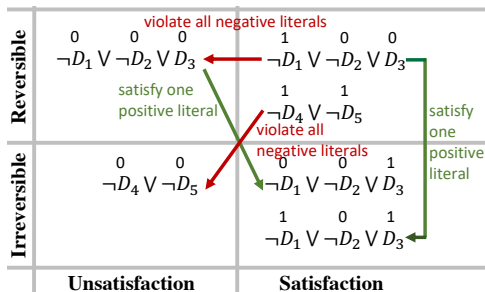
$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x}) \quad \text{where} \quad \sum_{i=1}^L C_i = L \quad (1)$$

By adding a high-cost penalty term for violated constraints, constrained optimization problem  $\rightarrow$  unconstrained problem:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x}) - \lambda' \sum_{i=1}^L (1 - C_i) \quad (2)$$

While exhaustive search is intractable, we use a beam-based search to find approximately optimal solutions for this objective.

# Constraint States



**Figure 3:** Clause states and possible transitions.  
 $D_i$  and  $\neg D_i$  denote positive and negative literal respectively

- S1: reversible unsatisfaction**
- S2: irreversible unsatisfaction**
- S3: reversible satisfaction**
- S4: irreversible satisfaction**

# Tracking Constraint States

## Prefix Tries

- Prefix trie,  $\mathcal{T}^+$  tracks *unsatisfied positive* literals from all clauses in states S1 and S3.
- Prefix trie,  $\mathcal{T}^-$  tracks *satisfied negative* literals from all clauses in state S3.

## How Prefix Tries Changes

- a positive literal satisfied  $\rightarrow$  its clause in state S1 or S3 henceforth irreversibly satisfied (state S4)  $\rightarrow$  remove all literals of that clause from both tries and stop tracking.
- a negative literal violated  $\rightarrow$  remove it from the trie  $\mathcal{T}^- \rightarrow$  switch back to S1 or S2 once all negative literals of a clause in state S3 has been removed.

# Algorithm

## High-level Intuition

At each time step, NEUROLOGIC selects generation hypotheses in consideration of both the objective function and the diversity of the partially satisfied constraints. We achieve such by 3 steps: *pruning*, *grouping*, and *selecting*.

**Pruning Step:** We first discard any  $h$  with irreversible unsatisfied clause (state S2) to focus only on candidates that might satisfy all constraints.

**Grouping Step:** Next, we select the beam from the pruned candidates.

**Selecting Step:** To select best ones from each group, we first rank candidates within a group by score function:

$$s = P_{\theta}(\mathbf{y}_t \mid \mathbf{y}_{<t}) + \lambda \cdot \max_{\substack{D(\mathbf{a}_i, \mathbf{y}) \\ \in \text{stateS1}}} \frac{|\hat{\mathbf{a}}_i|}{|\mathbf{a}_i|} \quad (3)$$

# Algorithm

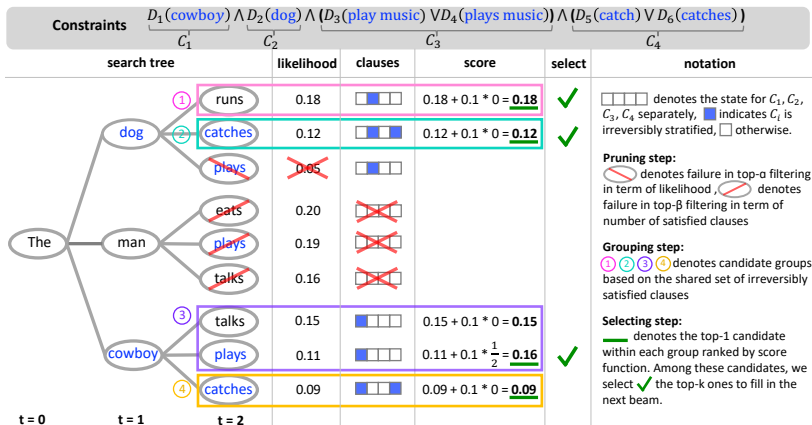


Figure 4: Illustration of the NEUROLOGIC decoding procedure. In this example,  $k = 3$ ,  $\alpha = 8$ ,  $\beta = 2$ ,  $\lambda = 0.1$ .

# Experiment I: Constrained Generation

## COMMONGEN [2]

COMMONGEN is a benchmark dataset designed as a test of generative commonsense reasoning. Given a set of common concepts (e.g., dog, frisbee, catch, throw); the task is to generate a coherent sentence describing an everyday scenario using these concepts (e.g., “a man throws a frisbee and his dog catches it”).

## Problem Formulation

The input is an unordered set of  $n$  concepts  $\mathbf{x} = \{a_1, a_2, \dots, a_n\}$ , where each concept  $a_i$  is a common object (noun) or action (verb). The expected output is a simple, grammatical sentence  $\mathbf{y} \in \mathcal{Y}$  that describes a common scenario using all given concepts in  $\mathbf{x}$  with correct morphological inflections.

[2] Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, Xiang Ren. CommonGen: A Constrained Text Generation Challenge for Generative Commonsense Reasoning. EMNLP2020, Findings.

# Results I: NEUROLOGIC vs Other Decoding Methods

Decode Method	ROUGE-L	BLEU-3/4	METEOR	CIDEr	SPICE	Coverage
Greedy Decoding	35.3	25.2 16.7	25.8	10.2	24.4	80.3
Top-k Sampling	33.8	22.5 14.4	24.9	9.2	22.7	79.4
Top-p Sampling	35.3	25.0 16.5	25.7	10.2	24.1	80.1
Beam Search	<u>40.3</u>	<u>34.2</u> <u>24.7</u>	<u>27.6</u>	<u>13.4</u>	<u>27.1</u>	82.2
Hokamp and Liu	37.6	25.6 16.8	25.9	11.1	25.1	97.2
Post and Vilar	38.3	28.1 18.6	26.7	11.8	26.0	<u>97.4</u>
Hu et al.	38.2	27.8 18.4	26.7	11.7	26.1	<u>97.4</u>
NEUROLOGIC	<b>42.8</b>	<b>36.7 26.7</b>	<b>30.2</b>	<b>14.7</b>	<b>30.3</b>	<b>97.7</b>

Table 1: Performance of different decoding methods using supervised GPT2-L on the COMMONGEN test set.

# Results II: NEUROLOGIC across Different Supervised Models

Model	ROUGE - L	BLEU - 3 & 4		METEOR	CIDEr	SPICE	Coverage
GPT-2	40.3 → 42.8	34.2 → 36.7	24.7 → 26.7	27.6 → 30.2	13.4 → 14.7	27.1 → 30.3	82.2 → 97.7
BERT-Gen	42.4 → 43.8	37.5 → 38.9	27.0 → 28.2	29.5 → 30.9	14.9 → 15.5	29.8 → <u>31.4</u>	89.2 → 97.3
UniLM	44.3 → <b>45.8</b>	40.6 → <b>42.8</b>	29.9 → <b>31.5</b>	30.1 → <b>31.7</b>	15.5 → <u>16.6</u>	30.6 → <b>32.5</b>	90.5 → 97.8
UniLM-v2	43.5 → 44.2	39.2 → 39.5	28.3 → 28.5	30.6 → <u>31.3</u>	15.2 → <b>16.8</b>	30.8 → 31.1	92.8 → 97.9
BART	43.3 → 44.7	39.9 → <u>41.3</u>	29.1 → <u>30.6</u>	30.4 → 31.0	15.2 → 15.9	30.6 → 31.0	95.0 → <b>98.7</b>
T5-Large	43.9 → <u>44.8</u>	36.6 → 38.5	26.9 → 28.1	28.9 → 30.7	14.3 → 15.5	29.5 → 30.8	89.7 → <u>98.5</u>

Table 2: Experimental results of different supervised models on the COMMONGEN test set.



# Results III: NEUROLOGIC with Unsupervised Models

Domain Adaption	Model	ROUGE - L	BLEU - 3 & 4		METEOR	CIDEr	SPICE	Coverage
No	GPT	26.7 $\rightarrow$ 41.3	3.0 $\rightarrow$ 25.1	1.1 $\rightarrow$ 15.9	9.2 $\rightarrow$ 28.8	0.9 $\rightarrow$ 11.7	8.0 $\rightarrow$ 29.7	8.4 $\rightarrow$ <b>97.4</b>
	GPT-2	19.7 $\rightarrow$ <b>42.9</b>	4.1 $\rightarrow$ <u>34.4</u>	1.5 $\rightarrow$ <u>23.5</u>	11.2 $\rightarrow$ <u>30.7</u>	0.4 $\rightarrow$ <u>13.6</u>	7.1 $\rightarrow$ <u>31.4</u>	8.3 $\rightarrow$ 96.0
Yes	GPT-2	29.8 $\rightarrow$ <u>42.4</u>	9.5 $\rightarrow$ <b>36.1</b>	4.0 $\rightarrow$ <b>25.1</b>	11.7 $\rightarrow$ <b>31.3</b>	1.7 $\rightarrow$ <b>13.9</b>	8.0 $\rightarrow$ <b>31.8</b>	9.3 $\rightarrow$ <u>96.1</u>

**Table 3:** Experimental results in zero-shot (unsupervised) setting on the COMMONGEN test set with and without language domain adaption.

# Results IV: Ablation

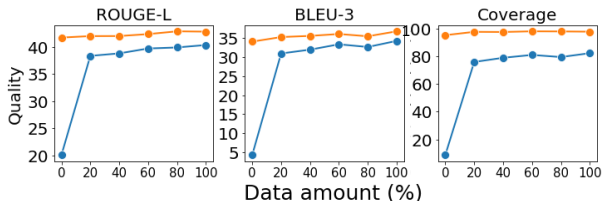


Figure 5: Performance (y-axis) of supervised GPT2-L on COMMONGEN, with a varying amount of training data for supervision (x-axis).

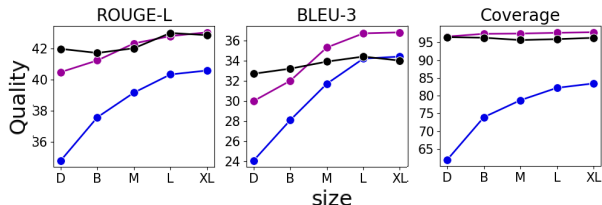


Figure 6: Performance (y-axis) of GPT-2 with varying model sizes (x-axis).

# Experiment II Results: Recipe Generation

Decode Method	ROUGE-L	BLEU-3/4	METEOR	Coverage	Extra
Top-k Sampling	27.5	15.2 9.5	19.2	84.8	16.0
Top-p Sampling	28.7	<u>17.6</u> 11.7	19.4	86.4	15.4
Beam Search	<u>29.4</u>	17.4 <u>12.0</u>	<u>19.7</u>	86.5	14.3
Post and Vilar	26.1	13.6 8.8	16.5	<u>89.6</u>	1.15
Hu et al	26.1	13.6 8.8	16.5	<u>89.6</u>	<u>1.13</u>
NEUROLOGIC	<b>32.1</b>	<b>19.5 13.8</b>	<b>19.8</b>	<b>95.8</b>	<b>0.6</b>

**Table 4:** Experimental results of different decoding methods with RecipeGPT on the Recipe1M+ test set. Coverage indicates the average percentage of ingredients that are covered in the generated recipe, while Extra corresponds to the average ratio of hallucinated ingredients over the number of given ingredients.

# An Example Illustration

Concept-Set {lose, board, balance, fall, ride}

Supervised Setting

## Decode with Beam Search

[GPT-2]: Someone loses balance and falls off his bike.

[UniLM]: A man is trying to keep his balance as he falls off a board.

[BART]: A man loses his balance and falls off the balance while riding a skateboard.

[T5]: a man loses his balance on the board and falls.

## Decode with NEUROLOGIC

[GPT-2]: A man loses his balance as he rides a roller coaster and falls off the board.

[UniLM]: Someone loses balance on the ride and falls off the balance board.

[BART]: A man loses his balance on a ride and falls off the board.

[T5]: a rider loses his balance and falls off the board.

Zero Shot Setting

## Decode with NEUROLOGIC

[GPT]: a woman lost her balance riding a horse, falling off the horse, and hitting her head on a board

[GPT-2]: The boy lost his balance riding the bike, falling off the bike, and hitting his head on the board.

Figure 7: Generated texts for the given concept-set.

# Contents

## ① NEUROLOGIC DECODING

## ② TEMPORAL REASONING

## ③ (COMET-)ATOMIC<sub>20</sub><sup>20</sup>

## ④ Hashtags, Emotions, and Comments

# TEMPORAL REASONING

## Temporal Reasoning on Implicit Events from Distant Supervision

**Ben Zhou<sup>1,2</sup> Kyle Richardson<sup>1</sup> Qiang Ning<sup>3</sup> Tushar Khot<sup>1</sup> Ashish Sabharwal<sup>1</sup> Dan Roth<sup>2</sup>**

<sup>\*1</sup>Allen Institute for AI <sup>2</sup>University of Pennsylvania <sup>3</sup>Amazon

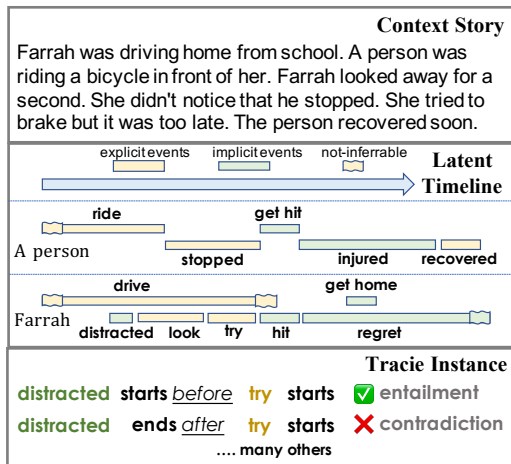
{kyler,tushark,ashishs}@allenai.org {xyzhou,danroth}@cis.upenn.edu qning@amazon.com

Figure 8: Paper information. [3]

---

[3] Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, D. Roth. Temporal Reasoning on Implicit Events from Distant Supervision. NAACL 2021.

# Task Definition



**Figure 9:** The task focuses on temporal relations on implicit events in short stories. A story, its latent timeline, and example TRACIE instances from it. For simplicity, events are shortened to single verbs and the timeline is exaggerated.

# Contributions

- A temporal relation dataset TRACIE focusing on **implicit events**.
- A **distant supervision process** for temporal understanding of implicit events.
- A **reasoning** model that makes end-time comparisons using predictions of start-time distances and durations.



# The TRACIE Dataset

Context Story (Premise)	Hypothesis	Inference Label
Tom needed to get braces. He was afraid of them. The dentist assured him everything would be fine. Tom had them on for a while. Once removed he felt it was worth it.	Tom avoids foods he can't eat with braces <b>starts</b> <b>before</b> the braces are removed.	entailment
We were all watching Spongebob as a family. It is a kid's show but all really enjoyed it. This one episode was especially funny for the adults. It has humor in it that is funny for kids and adults. It is something we can all watch...	The adults laughed at the jokes <b>ends</b> <b>before</b> we watch Spongebob as a family	contradiction
I was throwing the baseball with my son. He threw one past me that landed in the lake. I reached in to get the ball. I lost my balance and fell in. I got the ball and a bath all in one shot!	The ball was in the boys hand <b>starts</b> <b>after</b> he reached for the ball	contradiction

**Figure 10:** Example TRACIE instances. The **comparator**  $l \in \{\text{starts}, \text{ends}\}$  and **relation**  $r \in \{\text{before}, \text{after}\}$  in each hypothesis are highlighted, in addition to the corresponding explicit event from the story.

# The TRACIE Dataset




Illustration	Allen's Relation	Tracie's Relation
	Precedes, Meets	Starts Before Ends Before
	Overlaps, Finished-by, Contains, Starts, Equals, Started-by	Starts Before Ends After
	During, Finishes, Overlapped-by, Met-by, Preceded-by	Starts After Ends After

Figure 11: TRACIE's label definition and its relation to Allen's interval algebra, with a graph illustration between an **implicit event** and an **explicit event**.

# Implicit Event Generation

- 1 We randomly sample short stories from the ROCStories dataset [4].
- 2 For each story, one annotator writes 5 implicit event phrases that are not explicitly mentioned by the given story, but are inferable and relevant.
- 3 The annotator additionally rewrites two explicit events closest to the implicit event's start and end time, respectively.
- 4 With these two events, we can build two TRACIE instances (minus the *temporal-relation*) per implicit event, which accounts for 10 instances in total per story.

---

[4] Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. A Corpus and Evaluation Framework for Deeper Understanding of Commonsense Stories. NAACL2016.

# Pattern-Based Pre-Training

## High-level Intuition

we believe that it is more efficient to build a model that learns the prior knowledge needed for the task with distant signals and only subsequently learns the task definition through a small training set.

# Distant Supervision Collection

text

I went to the park on January 1<sup>st</sup>. I was very hungry after some hiking. Luckily, I purchased a lot of food before I went to the park. I enjoyed the trip and wrote an online review about the trip on the 10<sup>th</sup>.

within-sentence

[I purchased food, I went to the park.]: before

cross-sentence

[I went to the park, I wrote a review]: before, weeks

Figure 12: Extraction for start-time comparisons applied to an example paragraph.

We describe the sources of distant supervision signals with the goal of understanding the relative order between two events' start times as well as the relative distance between them.

# Supervision Instances Construction

- 1 Each instance comprises an event pair, a temporal relation, and an estimation on the temporal difference between the two start times.
- 2 Each event is a phrase constructed by taking all relevant arguments of the predicate verb in the SRL parses.
- 3 We represent the differences between the two start times as one of seven coarse temporal units:  $\{\leq \text{minutes}, \text{hours}, \text{days}, \text{weeks}, \text{months}, \text{years}, \geq \text{decades}\}$ .
- 4 In addition to the event pairs, we randomly sample sentences within the paragraph to use as the context that better defines the events.

# Pattern-Based Temporal Model (PTNTIME)

## Data Format

Input sequences `event:[EventA]`  
`starts[Relation][EventB].story:[Paragraph].`

Output sequences `answer:[Label][Distance].`

## PTNTIME

- We use a pre-trained sequence-to-sequence model as our base model and additionally pre-train this model using the data collected.
- PTNTIME serves as new set of *temporally-aware* model weights that can be used in place of existing pre-trained models and fine-tuned on TRACIE.

# TRACIE Task Formulation

comparator $l$	relation $r_l(e_1, e_2) =$
ends	before if $\text{end}_1 < \text{start}_2$
	after otherwise
starts	before if $\text{start}_1 < \text{start}_2$
	after otherwise

Figure 13: Decomposition of the relation functions that solve TRACIE instances (equal timepoints ignored).



# Neural-symbolic Model

## Two Modules

- Distance function:  $\text{dist}(e_i, e_j) = \mathbf{start}_i - \mathbf{start}_j$ .
- Duration function:  $\text{dur}(e_j) = \mathbf{duration}_j$ .

By exploiting the rule that an end point  $\mathbf{end}_j$  can be computed as  $\mathbf{end}_j = \mathbf{start}_j + \mathbf{duration}_j$ , we can, for example, decompose the relation  $r_{\text{ends}}(e_1, e_2) = \mathbf{before}$  (i.e.,  $e_1$  ends before  $e_2$ ) in terms of our two modules as follows via simple algebraic manipulation:

$$\begin{aligned}
 r_{\text{ends}}(e_1, e_2) &= \mathbf{before} \\
 &\Leftrightarrow \mathbf{end}_1 < \mathbf{start}_2 \\
 &\Leftrightarrow \mathbf{start}_1 + \mathbf{duration}_1 < \mathbf{start}_2 \\
 &\Leftrightarrow (\mathbf{start}_1 - \mathbf{start}_2) + \mathbf{duration}_1 < 0 \\
 &\Leftrightarrow \text{dist}(e_1, e_2) + \text{dur}(e_1) < 0
 \end{aligned}$$

# Distance Estimation

- We use the output from PTNTIME to approximate the function  $\text{dist}(\cdot)$ .
- Following the sequence formulation of PTNTIME, we replace [EventA] with the textual description of  $e_1$ , [EventB] with the textual description of  $e_2$ , and [Paragraph] with the context (premise), and fix [Relation] to be *before*. By taking the values of the vocabulary indices corresponding to “positive” and “negative” from the logits of [Label] and applying a softmax operation, we get  $P_{\text{before}}$  and  $P_{\text{after}}$ . These are the probability of  $e_1$  starting before and after  $e_2$ , respectively, and are used to define the vector  $\mathbf{p} = [P_{\text{before}}, P_{\text{after}}]$ .
- Similarly, we apply softmax to the logits of [Distance] over the 7 words representing the temporal units to obtain 7 values that approximate the probabilities of the distance between two events’ start times being closest to each temporal unit. We place the 7 values in temporal units’ increasing order in vector  $\mathbf{d}$ . To represent  $|\text{start}_1 - \text{start}_2|$  with a single value, we dot product the probabilities with an incremental constant vector  $\mathbf{c} = [0, 1, 2, 3, 4, 5, 6]$ . To get the direction, we apply the tanh function to the difference between the probabilities in  $\mathbf{p}$ .

# Duration Estimation

- To obtain a model to estimate  $\text{dur}(\cdot)$ , we pre-train a sequence-to-sequence model with the duration data from [5].
- The data contains over 1 million events with their corresponding duration values.
- We map each instance to an input sequence `event : [Event]` `story : [Story]` and a corresponding output sequence `answer : [Value]`.

---

[5] Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. Temporal Common Sense Acquisition with Minimal Supervision. ACL2020.

# Computation and Learning

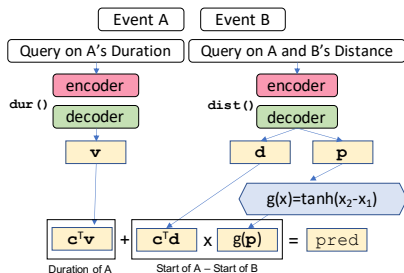


Figure 14: A schematic overview of SYMTIME to compare event  $A$ 's end time with event  $B$ 's start time via modular predictions about  $A$ 's duration and distance from  $B$  and their symbolic combination (bottom).

$$\begin{aligned} \text{dist}(\cdot) &= \text{start}_1 - \text{start}_2 \\ &= \mathbf{c}^T \mathbf{d} * \tanh(\text{INT}_{max} * (\mathbf{p}_2 - \mathbf{p}_1)) \end{aligned} \quad (4)$$

$$\text{dur}(\cdot) = \text{duration}_1 = \mathbf{c}^T \mathbf{v} \quad (5)$$

# Results I: I.I.D Setting

System	Start	End	All	Story
Majority	57.3	69.8	64.1	18.1
BiLSTM	53.7	63.5	59.1	10.9
Roberta-Large	78.5	78.3	78.4	26.1
T5-3B	79.4	77.4	78.3	26.9
BaseLM (T5-large)	75.5	75.4	75.4	22.6
BaseLM-MATRES	76.7	76.3	76.5	25.3
PTNTIME (ours)	81.4	77.5	79.3	31.0
SYMTIME (ours)	<b>82.1</b>	<b>79.4</b>	<b>80.6</b>	<b>32.0</b>
SYMTIME-ZEROSHOT	77.0	73.1	74.9	21.6

**Table 5:** Performance on TRACIE, best numbers in **bold**. BaseLM is T5-large; Story is the percentage of story-wide exact match; Majority is based on the comparator and temporal-relation distribution; Zeroshot uses no TRACIE instance as supervision.

# Results II: Uniform-prior Training Setting

System	Start	End	All	$\Delta$ All
Random	50.0	50.0	50.0	-14.1
BiLSTM	50.5	51.2	50.9	-8.2
Roberta-Large	75.1	68.1	71.3	-7.1
T5-3B	72.8	68.6	70.5	-7.8
BaseLM (T5-large)	68.1	67.8	67.9	-7.5
BaseLM-MATRES	76.3	69.9	72.8	-3.7
PTNTIME (ours)	80.6	73.2	76.6	-2.7
SYMTIME (ours)	<b>81.2</b>	<b>77.0</b>	<b>78.9</b>	-1.7
SYMTIME-ZEROSHOT	77.0	73.1	74.9	<b>0.0</b>

**Table 6:** Performance on TRACIEa uniform-prior training setting.  $\Delta$ All compares the difference with Table 5; Majority is equivalent to random guessing.

# Results III: MATRES for Explicit Events

System	OT-NS	OT	OT-MS	PT
Wang et.al.(2020)	85.9	-	-	-
BaseLM	86.0	87.5	77.4	69.0
PTNTIME	87.3	89.6	86.1	75.1

Table 7: Performance on MATRES[6] is not strictly comparable with the rest.

[6] Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. Joint constrained learning for event-event relation extraction. EMNLP2020.

# Results IV: Ablation

Sys.	BaseLM	PTN <b>TIME</b>	SYM <b>TIME</b>	Human
Acc.	52.6	72.2	75.3	82.5

Table 8: Performance on *no-story* TRACIE under the uniform-prior training setting.

Sys.	PTN <b>TIME</b>	cross-sentence	within-sentence
Acc.	80.6	79.9	63.7

Table 9: Comparison of pre-training data sources on TRACIE’s start time prediction accuracy, under the uniform-prior training setting.



# Contents

- ❶ NEUROLOGIC DECODING
- ❷ TEMPORAL REASONING
- ❸ (COMET-)ATOMIC<sub>20</sub><sup>20</sup>
- ❹ Hashtags, Emotions, and Comments

# (COMET-)ATOMIC<sub>20</sub>

## (COMET-)ATOMIC<sub>20</sub>: On Symbolic and Neural Commonsense Knowledge Graphs

Jena D. Hwang<sup>1\*</sup>, Chandra Bhagavatula<sup>1\*</sup>, Ronan Le Bras<sup>1</sup>, Jeff Da<sup>1</sup>, Keisuke Sakaguchi<sup>1</sup>,  
Antoine Bosselut<sup>13</sup> and Yejin Choi<sup>12</sup>

<sup>1</sup> Allen Institute for AI, WA, USA

<sup>2</sup> Paul G. Allen School of Computer Science & Engineering, WA, USA

<sup>3</sup> Stanford University, CA, USA

{jenah, chandrab, ronanl, jeffd, keisukes, antoineb, yejinc}@allenai.org

Figure 15: Paper information. [7]

[7] Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, Yejin Choi. COMET-ATOMIC 2020: On Symbolic and Neural Commonsense Knowledge Graphs. AAAI 2021.

# Motivation

- A new paradigm of language models as knowledge bases has emerged. In this setting, language models are **prompted** with natural language prefixes or questions, and they express knowledge through language generation.
- *Does scaling up language models actually endow them with commonsense knowledge?* They perform better when evaluated on knowledge bases that **prioritize ontological relations** and whose examples resemble language-like assertions (e.g., mango `ISA` fruit).
- Prior work has also shown that training language models on knowledge graph tuples leads them to learn to **express their implicit knowledge directly**, allowing them to provide commonsense knowledge on-demand.

# Contributions

- We present  $\text{ATOMIC}_{20}^{20}$ —**a new commonsense knowledge graph** covering social, physical, and eventive aspects of everyday inferential knowledge.
- We compare  $\text{ATOMIC}_{20}^{20}$  with other prominent CSKBs head-to-head and show that our new *symbolic* knowledge graph is **more accurate than any current CSKB**.
- We show that our new *neural* knowledge model COMET- $\text{ATOMIC}_{20}^{20}$  **successfully transfers  $\text{ATOMIC}_{20}^{20}$ 's declarative knowledge** to beat GPT-3, the largest pre-trained language model, in spite of using 400x fewer parameters.

# Comparisons of Three CSKGs

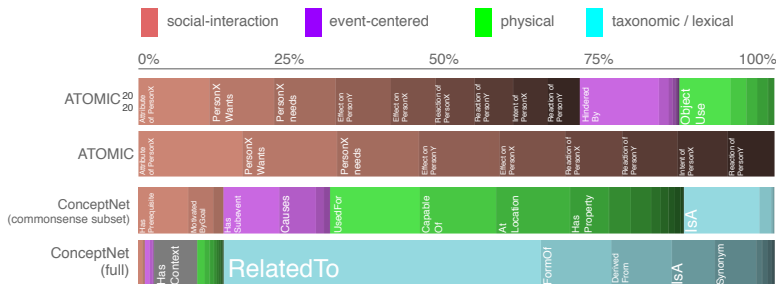


Figure 16: ATOMIC<sup>20</sup> tuple count distribution compared to ATOMIC and CONCEPTNET, either its commonsense subset or the full set.

# ATOMIC<sub>20</sub> Illustration

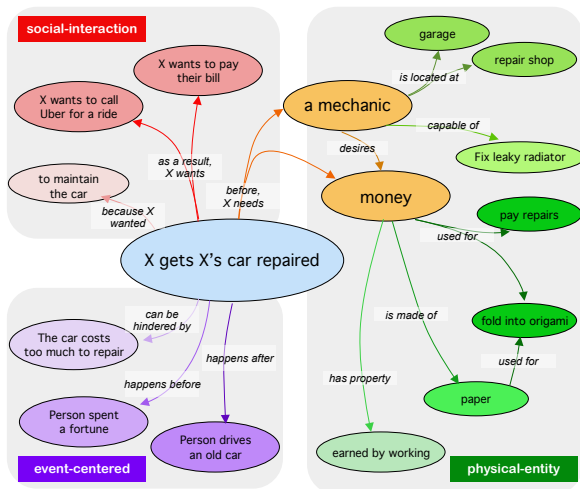


Figure 17: A tiny subset of ATOMIC<sub>20</sub>, a large atlas of social and physical commonsense relations. Relations in the top-left quadrant reflects relations from ATOMIC.

# ATOMIC<sub>20</sub> Relation

	Head	Relation	Tail	Size
PHYSICAL-ENTITY	bread	ObjectUse	make french toast	165,590
		AtLocation*	basket; pantry	20,221
		MadeUpOf	dough; wheat	3,345
		HasProperty*	cooked; nice to eat	5,617
	baker	CapableOf*	coat cake with icing	7,968
		Desires*	quality ingredients	2,737
		Not Desires*	bad yeast	2,838
EVENT-CENTERED	X runs out of steam	IsAfter	X exercises in the gym	22,453
		HasSubEvent	become tired	12,845
		IsBefore	X hits the showers	23,208
		HinderedBy	drinks too much coffee	106,658
		Causes	takes a break	376
		xReason	did not eat breakfast	334
	X runs out of steam	xNeed	do something tiring	128,955
		xAttr	old; lazy; lethargic	148,194
		xEffect	drinks some water	115,124
		xReact	tired	81,397
SOCIAL-INTERACTION	X votes for Y	xWant	to get some energy	135,360
		xIntent	to give support	72,677
		oEffect	receives praise	80,166
		oReact	grateful; confident	67,236
		oWant	thank X; celebrate	94,548

**Table 10:** Relations in ATOMIC<sub>20</sub> along with illustrative examples and their respective size. Relations that reflect semantically identical categories to CONCEPTNET is marked with an asterisk (\*).

# Accuracy Assessment

Knowledge Base	Accept	Reject	No Judgment
ATOMIC <sub>20</sub> <sup>20</sup>	<b>91.3</b>	<b>6.5</b>	2.2
ATOMIC	88.5	10.0	1.5
CONCEPTNET	88.6	7.5	3.9
TRANSOMCS	41.7	53.4	4.9

**Table 11: Accuracy** - Percentage (%) of tuples in the knowledge base evaluated by human crowdworkers as either always true or likely (Accept), farfetched/never or invalid (Reject), or unclear (No Judgment).



# Accuracy Assessment (Breakdown)

ATOMIC <sup>20</sup> <sub>20</sub>	ATOMIC	Relation	CN	T-OMCS
<b>92.3</b>		AtLocation*	89.4	34.3
<b>93.9</b>		CapableOf*	84.4	50.0
<b>94.6</b>		Causes	90.0	50.0
<b>96.9</b>		Desires*	96.3	48.2
<b>93.9</b>		HasProperty*	86.3	52.4
82.3		ObjUse/UsedFor	<b>96.3</b>	31.6
<b>98.5</b>		NotDesires*	96.3	
<b>96.9</b>		HasSubevent	88.1	57.7
		HasFirstSubevent	<b>93.8</b>	52.4
		HasLastSubevent	<b>95.6</b>	38.2
		HasPrerequisite	<b>94.4</b>	30.0
75.4		MadeUpOf/MadeOf	<b>88.1</b>	15.9
		PartOf	<b>71.9</b>	46.5
		HasA	<b>77.5</b>	43.5
<b>96.9</b>		HinderedBy		
<b>96.2</b>		isAfter		
<b>95.4</b>		isBefore		
<b>96.2</b>		isFilledBy		
		ReceiveAction	<b>84.4</b>	56.4
<b>91.5</b>	86.3	oEffect		
<b>91.5</b>	87.7	oReact		
88.5	<b>89.5</b>	oWant		
87.7	<b>91.0</b>	xAttr		
80.8	<b>87.2</b>	xEffect		
<b>93.1</b>	89.9	xIntent/MotivByGoal	84.4	27.1
<b>87.7</b>	85.1	xNeed		
90.8	<b>91.3</b>	xReact		
<b>96.2</b>		xReason		
82.3	88.4	xWant/CausesDesire	<b>90.0</b>	35.9

**Table 12:** KG accuracy values broken down by relation. Gray cells indicate statistically significant difference from ATOMIC<sup>20</sup><sub>20</sub> values. Dark gray cells signal instances where ATOMIC<sup>20</sup><sub>20</sub> values are significantly higher than its counterpart KB. Relational *cognates* have been grouped together and *exact matches* are asterisked (\*) (cf. Table 10).

# Coverage Assessment

Source KB↓	Target KB→			
	ATOMIC	CN	T-OMCS	ATOMIC <sub>20</sub> <sup>20</sup>
ATOMIC	-	0.1	0.0	100.0
CONCEPTNET	0.3	-	5.5	45.6
TRANSOMCS	0.0	0.4	-	0.3
ATOMIC <sub>20</sub> <sup>20</sup>	60.2	9.3	1.4	-

**Table 13: Coverage Precision** - Average number of times (in %) a tuple in Source KB is found in Target KB.

Source KB↓	Target KB→			
	ATOMIC	CN	T-OMCS	ATOMIC <sub>20</sub> <sup>20</sup>
ATOMIC	-	0.3	0.0	60.1
CONCEPTNET	0.1	-	0.3	8.9
TRANSOMCS	0.0	7.6	-	1.3
ATOMIC <sub>20</sub> <sup>20</sup>	100.1 <sup>†</sup>	47.8	0.4	-

**Table 14: Coverage Recall** - Average number of times (in %) a tuple in Target KB is found in Source KB.

<sup>†</sup>This value is greater than 100 because multiple tuples in ATOMIC<sub>20</sub><sup>20</sup> can map to the same tuple in ATOMIC.

# Human Evaluation (Crowdsource, \$15 per hour)

KG	Model	Accept	Reject	No Judgm.
ATOMIC <sub>20</sub>	GPT2-XL	36.6	62.5	0.9
	GPT-3	73.0	24.6	2.5
	COMET(GPT2-XL)	72.5	26.6	0.9
	COMET(BART)	<b>84.5</b>	<b>13.8</b>	1.7
ATOMIC	GPT2-XL	38.3	61.2	0.4
	COMET(GPT2-XL)	64.1	34.7	1.2
	COMET(BART)	<b>83.1</b>	<b>15.3</b>	1.6
CONCEPTNET	GPT2-XL	50.3	42.1	7.7
	COMET(GPT2-XL)	74.5	19.0	6.4
	COMET(BART)	<b>75.5</b>	<b>17.9</b>	6.6
TRANSOMCS	GPT2-XL	<b>28.7</b>	<b>53.5</b>	17.8
	COMET(GPT2-XL)	26.9	60.9	12.2
	COMET(BART)	23.8	65.9	10.3

**Table 15:** Human evaluation of generation accuracy (%). Each model uses greedy decoding to generate the *tail* of 5K randomly-sampled test prefixes (*head, relation*) from each knowledge graph. GPT2-XL, GPT-3 and BART have 1.5B, 175B and 440M parameters, respectively.

# Automatic Evaluation

		Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	ROUGE-L	CIDEr	BERT Score
ATOMIC <sub>20</sub>	COMET(GPT2-XL)	0.401	0.247	0.168	0.123	0.288	0.473	0.620	0.632
	COMET(BART)	<b>0.462</b>	<b>0.280</b>	<b>0.182</b>	<b>0.124</b>	<b>0.325</b>	<b>0.486</b>	<b>0.632</b>	<b>0.636</b>
ATOMIC	COMET(GPT2-XL)	0.429	0.300	0.225	<b>0.187</b>	0.297	0.527	0.754	0.638
	COMET(BART)	<b>0.521</b>	<b>0.330</b>	0.225	0.164	<b>0.351</b>	<b>0.552</b>	<b>0.766</b>	<b>0.650</b>
CONCEPTNET	COMET(GPT2-XL)	0.152	<b>0.115</b>	<b>0.092</b>	<b>0.080</b>	<b>0.131</b>	<b>0.193</b>	<b>0.421</b>	<b>0.552</b>
	COMET(BART)	<b>0.169</b>	0.108	0.069	0.046	0.127	0.180	0.350	0.532
TRANSOMCS	COMET(GPT2-XL)	0.298	0.000	0.000	0.000	0.179	0.300	0.249	0.677
	COMET(BART)	<b>0.351</b>	<b>0.216</b>	<b>0.004</b>	0.000	<b>0.201</b>	<b>0.352</b>	<b>0.298</b>	<b>0.681</b>

**Table 16:** Automated metrics for the quality of the *tail* generations for the knowledge models COMET(GPT2-XL) and COMET(BART). Each approach uses greedy decoding for all test prefixes for each KG. Similar results were obtained on the 5K sampled prefixes that were randomly selected for the human evaluation.

# Discussion I

## Do pretrained language models already encode commonsense knowledge?

The COMET training paradigm proposed by can perhaps be viewed less as a means of learning *knowledge* from KGs, and more as a method of learning an *interface* for language models to hypothesize encoded knowledge through language generation. We look forward to future work in this space that attempts to disentangle these two ideas.

# Discussion II

## What considerations should be made when designing commonsense knowledge resources?

Because certain types of knowledge are already encoded and expressible by pretrained language models, CSKG designers should focus on collecting examples and categories of knowledge that are less likely to be known by language models. For example, of the 378 test tuples evaluated by the GPT2-XL zero-shot model that contained the `HinderedBy` relation, only 1.3% were deemed plausible by human raters – jumping to 85% plausibility for COMET(BART) – pointing to an advantage in constructing  $\text{ATOMIC}_{20}^{20}$  with this relationship in mind or per-relation accuracy.

# Contents

## ① NEUROLOGIC DECODING

## ② TEMPORAL REASONING

## ③ (COMET-)ATOMIC<sub>20</sub><sup>20</sup>

## ④ Hashtags, Emotions, and Comments

# Hashtags, Emotions, and Comments

## Hashtags, Emotions, and Comments: A Large-Scale Dataset to Understand Fine-Grained Social Emotions to Online Topics

**Keyang Ding   Jing Li\*   Yuji Zhang**

Department of Computing, The Hong Kong Polytechnic University, HKSAR, China

`{keyang.ding, yu-ji.zhang}@connect.polyu.hk`

`jing-amelia.li@polyu.edu.hk`

Figure 18: Paper Info [8]

---

[8] Keyang Ding, Jing Li and Yuji Zhang. Hashtags, Emotions, and Comments: A Large-Scale Dataset to Understand Fine-Grained Social Emotions to Online Topics. EMNLP2020.



# Motivations

- Most of the related work focus on the feelings from writers and the existing studies concerning reader emotions mostly tackle well-written texts, such as news reports.
- Limited work has been done to characterize collective feelings from the public (henceforth **social emotions**) to an online topic described with fragmented and colloquial social media language.
- Where some previous efforts gather viewpoints from limited readers through user replies manual annotations, we focus on social emotions reflecting aggregated feelings from large amount of people.

# Task Definition

[H]:#张艺兴整蛊GAI#  
[T]: Lay played tricks on GAI.  
[E]: 😂: lol; 🤔: facepalm; 🐶: doge (tease).

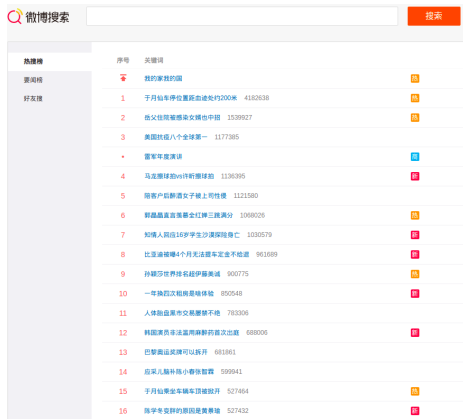
Figure 19: A Weibo hashtag and its resulting social emotions.

# Data Collection

Our dataset is built based on a Weibo emotion vote, where it provides users to vote for an emoji from a total of 24 emojis in the form of a questionnaire to represent their feelings to a trending hashtag.

# Data Collection

First, we tracked the trending hashtags following the everyday Weibo topic summary list from Apr to May 2020.



The screenshot shows the Weibo search interface with a search bar and a list of trending topics. The topics are ranked by popularity, with the top topic being '我的家我的国' (My Home, My Country).

热搜榜	序号	关键词	热度
热搜榜	1	我的家我的国	4182638
热搜榜	2	于月仙车停位置距血迹处约200米	1539927
热搜榜	3	岳父住院被感染女婿也中招	1177385
热搜榜	4	霍军年度演讲	1196395
热搜榜	5	马龙发球跑v8待新发球跑	1121580
热搜榜	6	随客户后醉酒女子被上男性侵	1069026
热搜榜	7	郭晶晶宣布霍金红牌三跳满分	1030579
热搜榜	8	知情人称16岁学生沙漠探险身亡	961689
热搜榜	9	比豆腐被曝4个月无法厘清定金不退还	900775
热搜榜	10	孙颖莎世界排名反超伊藤美诚	850548
热搜榜	11	一年换四次相房是啥体验	783306
热搜榜	12	人体胎盘黑市交易屡禁不绝	688006
热搜榜	13	韩国演员非法滥用麻醉药首次出庭	681861
热搜榜	14	巴黎奥运会奖牌可以拆开	599941
热搜榜	15	应采儿陈柏霖小春你智慧	527464
热搜榜	16	于月仙乘坐车辆车跟被脱开	527432

Figure 20: Illustration of the summary list.

# Data Collection

Then, we searched and parsed their emotion vote webpage via querying the hashtag in HTTP requests with the selenium package.



Figure 21: Illustration of the vote webpage.

# Data Collection

- Next, the crawled pages were parsed and analyzed using lxml package to gather the topics' emotion voting results. At last, hashtags with less than 100 voters were removed to filter out biased results.
- As Weibo only keeps emotions gaining the top three votes, we will hence focus on the top three emotions in the following discussions. These emotions were selected by over 83% voters on average and can still reflect feelings from the majority.
- Furthermore, to access the contexts of hashtags, we collected some user comments involved in a hashtag's discussion.

# Data Statistics

Dataset	Size	Len	Voters	Emos
Zhou et al. (2018)	5,586	702.4	157	6
Bostan et al. (2019)	5,000	11.3	331	8
Our dataset	13,766	5.4	3,250	24

Figure 22: Statistics: our data vs. prior resource. Size and emos are the number of instances and emotion types. Len and voters are the average number of words (after Chinese word segmentation) and the involved voters per instance.

# Data Statistics

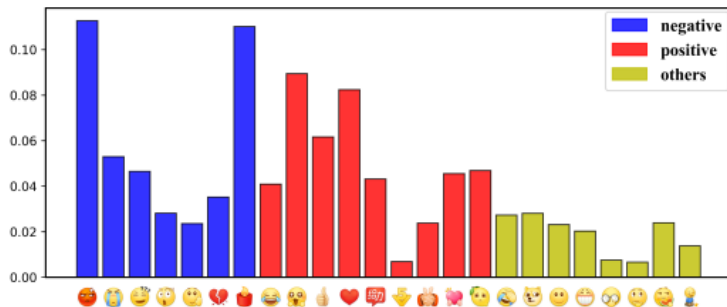
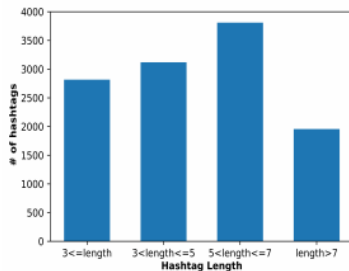
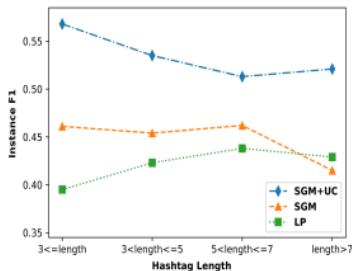


Figure 23: User preferences over varying emotions.



# Result Analyses



(a) Length vs. Instance F1      (b) Length vs. Hashtag Count

Figure 24: Instance F1 (left y-axis) in prediction and training hashtag number (right y-axis) over hashtag length (Chinese word count shown in x-axis).

Thanks for Listening.



NUSTM

<http://www.nustm.cn/member/rxia/index-cn.html>

<https://github.com/NUSTM>

# References I

- [1] X. Lu, P. West, R. Zellers, R. Le Bras, C. Bhagavatula, and Y. Choi, “NeuroLogic Decoding: (Un)supervised Neural Text Generation with Predicate Logic Constraints,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, (Online), pp. 4288–4299, Association for Computational Linguistics, June 2021.
- [2] B. Y. Lin, W. Zhou, M. Shen, P. Zhou, C. Bhagavatula, Y. Choi, and X. Ren, “CommonGen: A constrained text generation challenge for generative commonsense reasoning,” in *Findings of the Association for Computational Linguistics: EMNLP 2020*, (Online), pp. 1823–1840, Association for Computational Linguistics, Nov. 2020.
- [3] B. Zhou, K. Richardson, Q. Ning, T. Khot, A. Sabharwal, and D. Roth, “Temporal Reasoning on Implicit Events from Distant Supervision,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, (Online), pp. 1361–1371, Association for Computational Linguistics, June 2021.

# References II

- [4] N. Mostafazadeh, N. Chambers, X. He, D. Parikh, D. Batra, L. Vanderwende, P. Kohli, and J. Allen, “A Corpus and Evaluation Framework for Deeper Understanding of Commonsense Stories,” in *NAACL*, 2016.
- [5] B. Zhou, Q. Ning, D. Khashabi, and D. Roth, “Temporal Common Sense Acquisition with Minimal Supervision,” in *ACL*, 2020.
- [6] H. Wang, M. Chen, H. Zhang, and D. Roth, “Joint constrained learning for event-event relation extraction,” in *EMNLP*, 2020.
- [7] J. D. Hwang, C. Bhagavatula, R. L. Bras, J. Da, K. Sakaguchi, A. Bosselut, and Y. Choi, “(Comet-) Atomic 2020: On Symbolic and Neural Commonsense Knowledge Graphs,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 6384–6392, May 2021. Number: 7.
- [8] K. Ding, J. Li, and Y. Zhang, “Hashtags, Emotions, and Comments: A Large-Scale Dataset to Understand Fine-Grained Social Emotions to Online Topics,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, (Online), pp. 1376–1382, Association for Computational Linguistics, Nov. 2020.