





Information Sciences Institute

# Modeling Naive Psychology of Characters in Simple Commonsense Stories

#### Hannah Rashkin, Antoine Bosselut, Maarten Sap, Kevin Knight & Yejin Choi

Paul G. Allen School of Computer Science and Engineering, University of Washington Allen Institute for Artificial Intelligence Information Sciences Institute, University of Southern California

## Inferring Character State



## Reasoning about Naïve Psychology

New Story Commonsense Dataset:

- Open text + psychology theory
- Complete chains of mental states of characters
- Implied changes to characters
- Contextualized reasoning

C 🟠 Secure   https://uwnlp.github.io/storycommonsense/						
Browsing tool						
Select one of our stories and click on individual characters to see our annotations. For dev/test stories, hover over categories to see descriptions.						
Alicia's Tattoo 🔹						
Alicia loved tattoos. (Alida) Tattoo artist Alicia's motivation • to have a tattoo (Maslow: spiritual growth, love; Reiss: indep) Alicia's emotion • excited (joy, fear, sadness) • happy (joy, trust)						
She decided to get one. Alicia Tattoo artist						
She did research and found a great tattoo artist. Alicia Tattoo artist						
Alicia and the tattoo artist worked together to create a good design. Alicia Tattoo artist						
Alicia is now very happy with her new tattoo. Alicia Tattoo artist						

#### https://uwnlp.github.io/storycommonsense/

## How do we represent naïve psychology?

The band instructor told the band to start playing. <u>He often stopped the music when players were off-tone.</u>



## Naïve Psychology Annotations

- Motivation:
  - Causal source to actions
  - Motivational theories



- Emotional Reaction:
  - Causal effect of actions
  - Theories of emotion



### Motivation: Maslow Hierarchy of Needs (1943)



## Motivation: Reiss Categories (2004)



## Emotional Reaction: Plutchik (1980)

#### Plutchik's Wheel



Suddenly, they heard a loud noise. ↓ feel Fear, Surprise

## Implicit Mental State Changes

The band instructor told the band to start playing.

He often stopped the music when players were off-tone.

They grew tired and started playing worse after a while.

The instructor was furious and threw his chair.

How are players affected?

- $\rightarrow$  implicitly involved
- $\rightarrow$  inference in these cases

## **Tracking Mental States**

The band instructor told the band to start playing.

He often stopped the music when players were off-tone.

They grew tired and started playing worse after a while.

The instructor was furious and threw his chair.

He cancelled practice and expected us to perform tomorrow.

Why does the instructor cancel practice?

- $\rightarrow$  based on previous info
- $\rightarrow$  need to incorporate context

## **Related Work**

- Reasoning about narratives (Mostafazadeh et al 2016)
- Detecting emotional content (Mohammad et al 2013) or stimuli (Gui et al 2017) of a statement

Our work:

- Both motivation and emotion for a character's outlook
- Leverage psychology theories and natural language explanations

## Full Annotation Chain





## Full Annotation Chain







## Full Annotation Chain

Split into multiple stages



### Character Identification



#### Motivation



#### **Emotional Reaction**



## Data Collection Summary

Over 300k low-level annotations for 15k stories from ROC training set

		Open-text	Open-text + categories		
		train	dev	test	
# cl	haracter-line pairs	200k	25k	23k	
And the second s	w/ motivation change	40k	9k	7k	>50k motiv. changes
	w/ emotional reaction change	77k	15k	14k	>100k emotion changes

# Annotated Data Distributions (Motivation)



- Fair amount of diversity in the open-text
- ~1/3 have positive motivation change:



## Annotated Data Distributions (Emotion)



• Lots of happy stories

Sampled Open-text

• ~2/3 have positive emotion change:





% Annotations where selected

New Tasks

Given a story excerpt and a character can we explain the mental state:

- <u>Explanation Generation</u>: Generate open-text explanation of motivation/emotional reaction
- <u>State Classification</u>: Predict Maslow/Reiss/Plutchik category

## Task 1 - Explanation Generation

Explain mental state of character using natural language



## Modeling

- Using encoder-decoder framework
- Encoders LSTM, CNN, REN, NPN
- Decoder for generation: single layer LSTM



## Encoding Modules

Given entity  $e_j$  and line  $x^s$  (and entity-specific context sentences  $x^c[e_j]$ )  $h = f_{enc}(x^s, x^c[e_j])$ 

Encoding functions:

 CNN, LSTM: encode last line and context -- concatenate

## **Entity Modeling**

- Recurrent Entity Networks (Henaff et al 2017)
  - Store separate memory cells for each story character
  - Update after each sentence with sentence-based hidden states

- Neural Process Networks (Bosselut et al 2018)
  - Also has separate representations for each character
  - Updates after each sentence using learned action embeddings

## Explanation Generation Set-up

Evaluation: Cosine similarity of generated response to reference

Random baseline: Select random answer from dev set

- Responses are short/formulaic
- Words for describing intent/emotion are close in embedding space

#### **Explanation Generation Results**

Cos. Similarity to Reference



### Task 2 – Mental State Classification

Predicting psychological categories for mental state



## Modeling

- Using encoder-decoder framework
- Encoders LSTM, CNN, REN, NPN
- Decoder for categorization:
  logistic regression



## State Classification Set-Up



- 80% of dev set tuning predictions
- Each category as binary variable
- F1 taking # true positives across all classes

$$Recall = \frac{\# True Positive}{\# Actual Positive} Precision = \frac{\# True Positive}{\# Predicted Positive}$$

$$F_1 = \frac{2}{\frac{1}{\text{prec}} + \frac{1}{\text{rec}}}$$

## State Classification Results

- CNN and LSTM perform best on motivation categories
- Entity modeling has slight improvement in Plutchik



#### F1 Performance

#### Further Improvement

#### F1 Performance



## Effect of Entity Specific Context

Including previous lines from context that include entity

Entity specific context: improves all models F1 by about 3-5%

#### F1 w/ and w/o context



## Pre-training Encoders

We have more open-text explanations than category annotations:

- 1. Pre-train encoders on opentext explanations
- 2. Fine-tune with the categorical labels



## Effect of Pretrained Encoders

F1 w/ and w/o Pretrained Encoders



## Performance Per Category

Highest performance:

- Frequent classes (eg. "joy" F1: 38.9%)
- Very concrete sets of actions ("physiological" F1: 40%)



### Future Work

- **Outside Knowledge**: Help with infrequent classes and subtle implied changes
- Social Commonsense: Help with inferring mental state especially in more contextual cases
- **Potential Applications**: Improving language models, chat systems, natural language understanding

## Conclusions

- New Dataset:
  - 15k roc stories annotated per character
    - >50k motivation changes
    - >100k emotions changes
  - o https://uwnlp.github.io/storycommonsense/

## Maslow Performance Per Class

- In general, more concrete, low-level categories are easier to predict
- Eg. Hunger is pretty easy ot identify even with surface level features



## Plutchik Performance Per Class

- Joy is easiest to distinguish
- More infrequent emotions like disgust is weaker performance

