

# Modeling the Continuum of Emotions in Neural Dialogue Systems

Nabiha Asghar

PhD Seminar  
19th February 2019

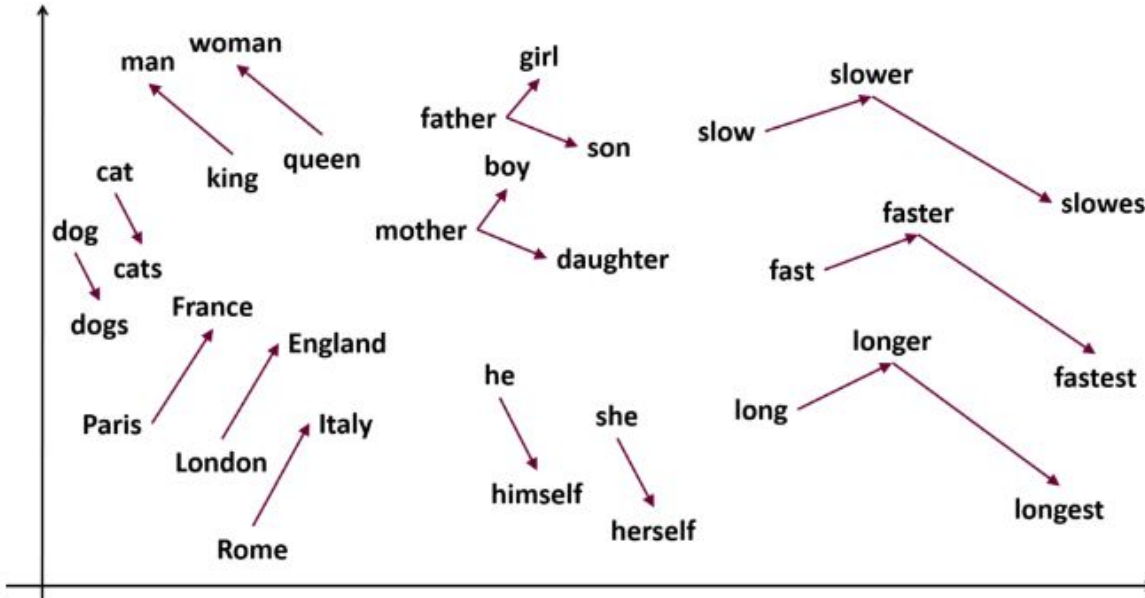
# Motivation

- Language → essential for human-human communication
- Humans are rational/logical *and emotional*
- Human actions/decisions are motivated by **both** emotional *and* non-emotional goals [1,2] [Example: language use]
- Discrete Emotion Theory vs Continuous Emotion Theory
  - Most recent NN research studies discrete emotions

[1] R. Picard. "Affective Computing". MIT Press, 1997.

[2] Zhu & Thagard. "Emotion and Action". Philosophical Psychology Vol 15 No 1, 2002.

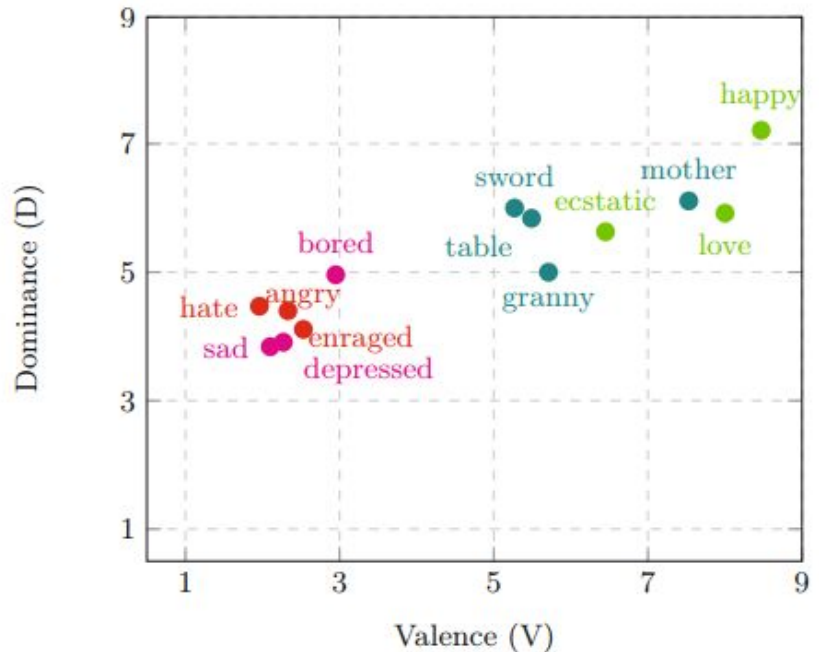
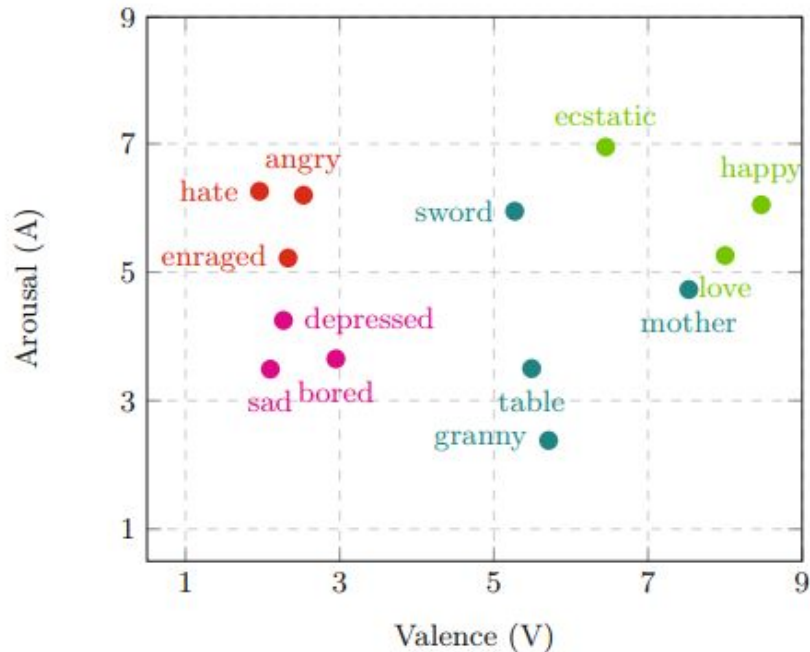
# Word Embeddings (Word2Vec)



Co-occurrence statistics insufficient to capture emotional features

Words different in sentiment often share context (e.g., “a good book” vs. “a bad book”).

# Affective Word Embeddings

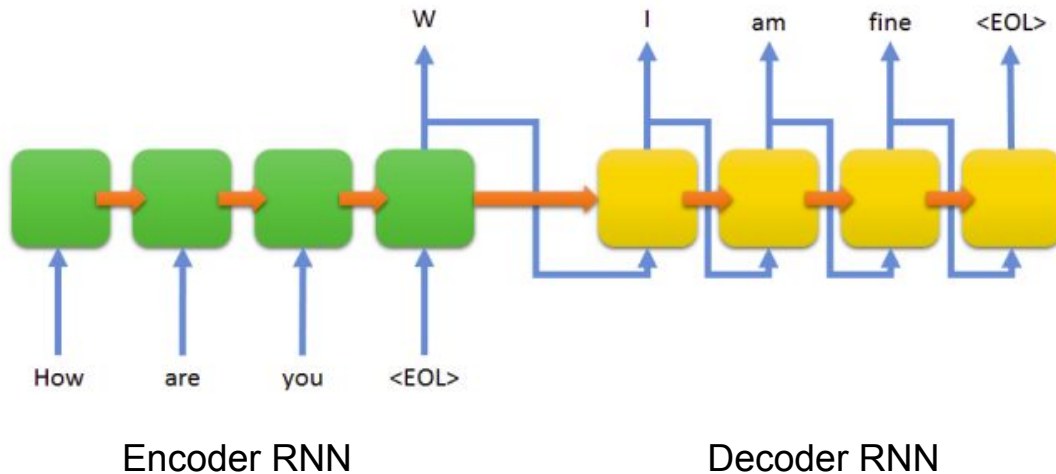


Valence -- Arousal -- Dominance  $\Rightarrow$  (VAD space)  $\Rightarrow$  Scale [0, 9]

# Goal

- Leverage word-level affect to generate emotional responses in dialogue
- Design:
  - Affective word embeddings
  - Affective loss functions
  - Affectively diverse “decoding” of response during inference
  - Use “Affect Control Theory” to generate emotional responses

# Recurrent NN Model (Seq2Seq)



Given a message response pair  $(X, Y)$  where  $X = x_1, \dots, x_m$  and  $Y = y_1, \dots, y_n$

$$L_{\text{XENT}}(\theta) = -\log p(Y|X) = -\sum_{i=1}^n \log p(y_i|y_1, \dots, y_{i-1}, X)$$

# Affective Word Embeddings

- Use cognitively engineered VAD dictionary [5]
- Define

$$W2AV(w) = \begin{cases} VAD(l(w)), & \text{if } l(w) \in dict \\ \boldsymbol{\eta} = [5, 1, 5], & \text{otherwise} \end{cases}$$

- Concatenate **W2AV** with **Word2Vec**
- Input the result to encoder and the decoder

# Affective Loss Functions

## Minimizing (Maximizing) Affective Dissonance

$$L_{\text{DMIN}}^i(\theta) = -(1 - \lambda) \log p(y_i | y_1, \dots, y_{i-1}, X) + \lambda \hat{p}(y_i) \left\| \sum_{j=1}^{|X|} \frac{W2AV(x_j)}{|X|} - \sum_{k=1}^i \frac{W2AV(y_k)}{i} \right\|_2$$

Average affect vector of the source sentence

Average affect vector of the target sub-sentence generated up to the current time step  $i$



# Affective Loss Functions

## Maximizing Affective Content

$$L_{AC}^i(\theta) = -(1 - \lambda) \log p(y_i | y_1, \dots, y_{i-1}, X) - \lambda \hat{p}(y_i) \|W2AV(y_i) - \boldsymbol{\eta}\|_2$$

These loss functions are not directly differentiable (because  $W2AV(x)$  is not continuous function, it's a dictionary. It is not a learnable function). So we relax it with predicted probability.

# Beam Search

- BS maintains top B most likely (sub)sequences
- At time t, augments the top-B subsequences from time t - 1 with all possible actions available
- Retain the top-B most likely branches up to time t (prune the rest)

$V \rightarrow$  vocabulary of tokens

$X \rightarrow$  input sequence

$\mathbf{y}_{i,[t-1]} \rightarrow$  i'th beam stored at time t-1

$Y_{[t-1]} = \{\mathbf{y}_{1,[t-1]}, \dots, \mathbf{y}_{B,[t-1]}\} \rightarrow$  set of beams stored at time t-1

$Y_{[t-1]} \times V \rightarrow$  set of all possible extensions of the beams from time t-1

$$Y_{[t]} = y_{1..t}^{1*}, \dots, y_{1..t}^{B*} = \underset{\substack{\mathbf{y}_{1,[t]}, \dots, \mathbf{y}_{B,[t]} \\ \in Y_{[t-1]} \times V}}{\arg \max}} \sum_{b=1}^B \sum_{i=1}^t \log p(y_{b,i} | \mathbf{y}_{b,[i-1]}, X)$$

# Diverse Beam Search

- Incorporate diversity among candidate outputs
- Divide top-B beams into G groups
- Measure the dissimilarity between group g and previous groups 1,  $\dots$ , g - 1 if token  $y_t$  is selected to extend any beam in group g

$$Y_{[t]}^g = \arg \max_{\substack{\mathbf{y}_{1,[t]}^g, \dots, \mathbf{y}_{B',[t]}^g \\ \in Y_{[t-1]}^g \times V}} \sum_{b=1}^{B'} \sum_{i=1}^t \log p(y_{b,i}^g | \mathbf{y}_{b,[i-1]}^g, X) + \lambda_g \Delta(Y_{[t]}^1, \dots, Y_{[t]}^{g-1})[y_{b,t}^g]$$

# Word-level Affective Diversity

$$\Delta_W(Y_{[t]}^1, \dots, Y_{[t]}^{g-1})[y_{b,t}^g] = - \sum_{j=1}^{g-1} \sum_{c=1}^{B'} \text{sim}(W2AV(y_{b,t}^g), W2AV(y_{c,t}^j))$$

$y_{b,t}^g$  → token under consideration at the current time step  $t$  for beam  $b$  in group  $g$

$y_{c,t}^j$  → token chosen for beam  $c$  in a previous group  $j$  at time  $t$

This metric ensures that the word affect at time  $t$  is diversified across groups

# Sentence-level Affective Diversity

$$\Delta_S(Y_{[t]}^1, \dots, Y_{[t]}^{g-1})[y_{b,t}^g] = - \sum_{j=1}^{g-1} \sum_{c=1}^{B'} \text{sim}(\Psi(\mathbf{y}_{b,[t]}^g), \Psi(\mathbf{y}_{c,[t]}^j))$$

where  $\Psi(\mathbf{y}_{i,[t]}^k) = \sum_{w \in \mathbf{y}_{i,[t]}^k} \text{W2AV}(w)$

- Computes the cumulative dissimilarity (given by the function  $\Psi$ ) between the current beam and all the previously generated beams in other groups
- Bag of affective words

# Experiments

- User Study (5 human judges rated responses)
  - Syntactic Coherence
  - Naturalness
  - Emotional Appropriateness
  - Syntactic Diversity
  - Affective Diversity

# Experiments

**Table 1.** The effect of affective word embeddings as input.

Model	Syntactic coherence	Natural	Emotional approp.
Word emb. (baseline)	1.48	0.69	0.41
Word+Affective emb.	<b>1.71</b> ↑	<b>1.05</b> ↑	<b>1.01</b> ↑

**Table 2.** The effect of affective loss functions.

Model	Syntactic coherence	Naturalness	Emotional approp
$L_{XENT}$ (baseline)	1.48	0.69	0.41
$L_{DMIN}$	<b>1.75</b> ↑	0.83 ↑	0.56 ↓
$L_{DMAX}$	1.74 ↑	0.85 ↑	0.58 ↑
$L_{AC}$	1.71 ↑	<b>0.95</b> ↑	<b>0.71</b> ↑

**Table 3.** Effect of affectively diverse decoding. H-DBS refers to Hamming-based DBS used in [22]. WL-ADBS and SL-ADBS are the proposed word-level and sentence-level affectively diverse beam search, respectively.

Model	Syntactic diversity	Affective diversity	# Emotionally approp. responses
BS (baseline)	1.23	0.87	0.89
H-DBS	1.47 ↑	0.79 ↓	0.78 ↓
WL-ADBS	<b>1.51</b> ↑	1.25 ↑	1.30 ↑
SL-ADBS	1.45 ↑	<b>1.31</b> ↑	<b>1.33</b> ↑

**Table 4.** Combining different affective strategies.

Model	Syntactic coherence	Naturalness	Emotional approp.
Traditional Seq2Seq (baseline)	1.48	0.69	0.41
Seq2Seq+Affective embeddings	1.71 ↑	1.05 ↑	1.01 ↑
Seq2Seq+Affective emb. & Loss	<b>1.76</b> ↓	1.03 ↓	1.07 ↑
Seq2Seq+Affective emb. & Loss & Decoding	1.69 ↓	<b>1.09</b> ↑	<b>1.10</b> ↓



# Better Heuristic?

- use Affect Control Theory to model emotional relationship between prompt and response
- Socio-mathematical theory of interaction between two human identities
- Example: friend-friend vs friend-enemy

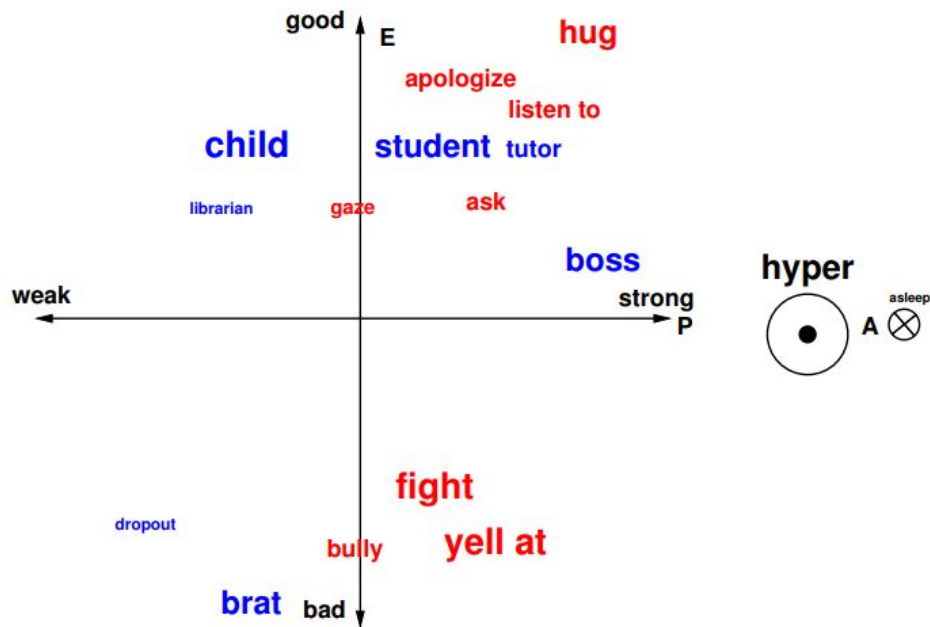
# Affect Control Theory (ACT)

- Emotions are points in 3D continuous space  $[-4.3, 4.3]^3$
- EPA space: **E**valuation (good/bad), **P**otency (strong/weak), **A**ctivity (excited/calm)
- Fundamental sentiments (**F**) vs Transient impressions (**T**)
- $\mathbf{f}(\text{mother}) = [2.9, 1.5, 0.6]$
- “A mother hugs a child”,  $\mathbf{t}(\text{mother}) = [3.5, 1.9, 0.85]$
- “A mother hits a child”,  $\mathbf{t}(\text{mother}) = [-1.0, 3.5, 2.2]$

# Affect Control Theory (ACT)

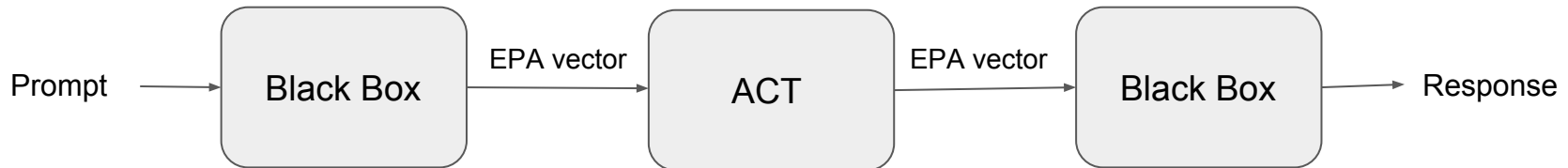
- Emotions are points in 3D continuous space  $[-4.3, 4.3]^3$
- EPA space: **Evaluation** (good/bad), **Potency** (strong/weak), **Activity** (excited/calm)
- Fundamental Sentiments **F**
- Transient Impressions **T**
- Deflection:  $D = \sum_i w_i (f_i - \tau_i)^2$

**Affect Control Principle:** actors work to minimize deflection, i.e. experience transient impressions that are consistent with their fundamental sentiments



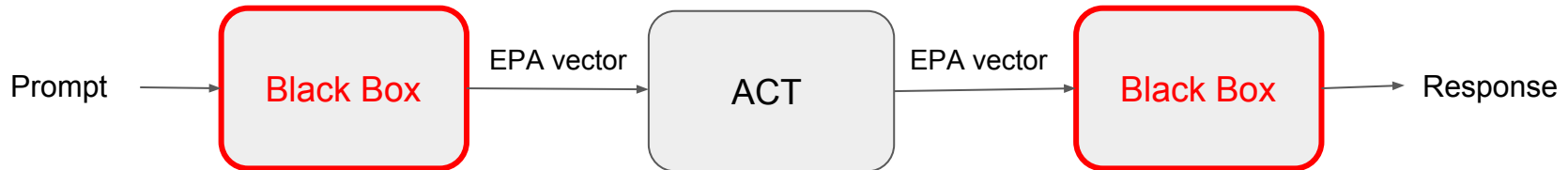
# Model

- ACT takes EPA vectors as input, produces EPA vectors (actions)
- Need a way to convert sentences into EPA, and vice versa
- Pipeline:



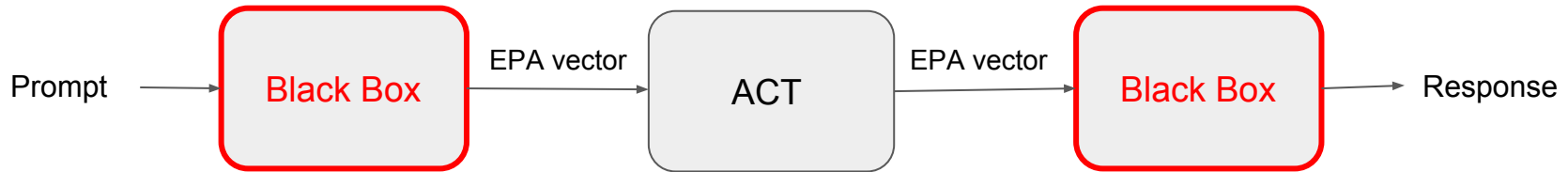
# Model

- ACT takes EPA vectors as input, produces EPA vectors (actions)
- Need a way to convert sentences into EPA, and vice versa
- Pipeline:



# Research Questions

1. How to convert prompt sentence into EPA vector?
2. How to convert EPA action to output sentence?

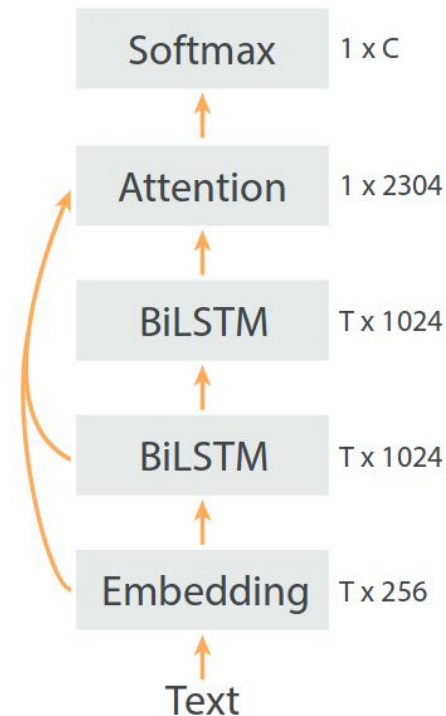


# Sentence to EPA

- We can train our own neural network, but labelling is too tedious

# Sentence to EPA

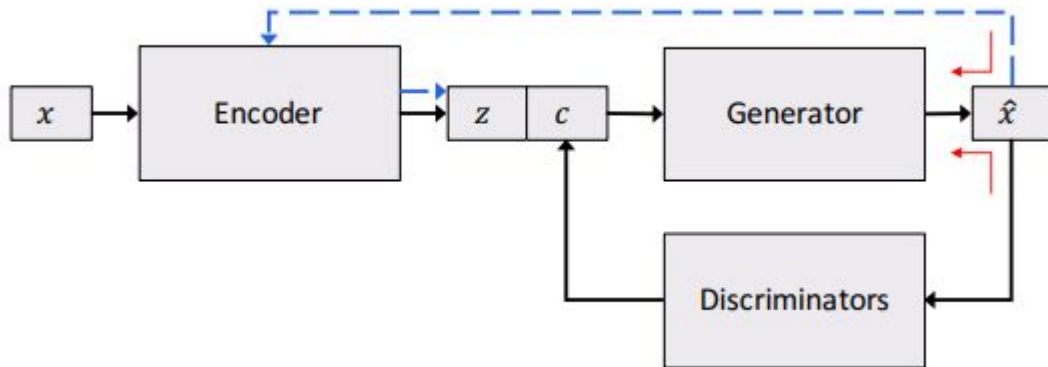
- We can train our own neural network, but labelling is too tedious
- Use **DeepMoji** - a pre-trained RNN model (billions of tweets). Predicts emojis (64 classes) given input sentence.
  - Asked 2 human judges to label the 64 emojis with EPA vectors. Average the labels.
  - This gives us an EPA for each emoji.
  - Given some input sentence, query DeepMoji and take weighted average of output



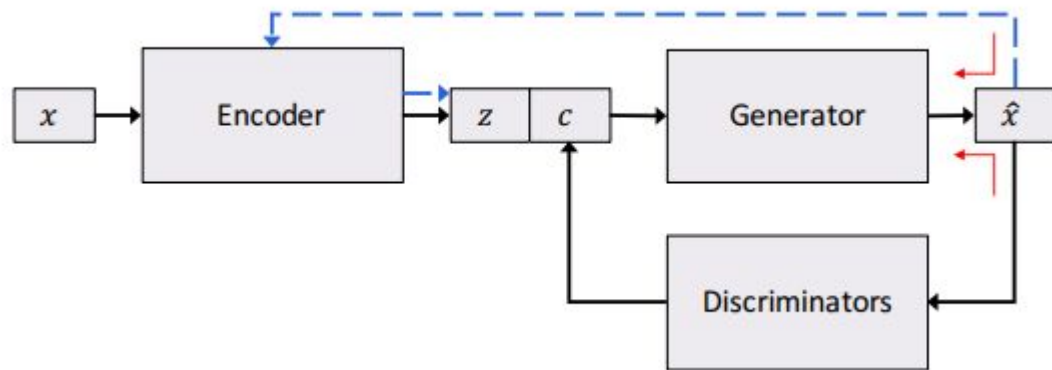


# EPA to Sentence

- Variational AutoEncoder (VAE) and discriminator to learn disentangled representation of emotions in text. Then query Generator with EPA



# Model



Generator Objective:

$$\min_{\theta_G} \mathcal{L}_G = \mathcal{L}_{\text{VAE}} + \lambda_c \mathcal{L}_{\text{Attr},c} + \lambda_z \mathcal{L}_{\text{Attr},z}$$

$$\mathcal{L}_{\text{VAE}}(\theta_G, \theta_E; \mathbf{x}) = \text{KL}(q_E(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) - \mathbb{E}_{q_E(\mathbf{z}|\mathbf{x})q_D(\mathbf{c}|\mathbf{x})} [\log p_G(\mathbf{x}|\mathbf{z}, \mathbf{c})]$$

$$\mathcal{L}_{\text{Attr},c}(\theta_G) = -\mathbb{E}_{p(\mathbf{z})p(\mathbf{c})} \left[ \log q_D(\mathbf{c}|\tilde{G}_\tau(\mathbf{z}, \mathbf{c})) \right]$$

$$\mathcal{L}_{\text{Attr},z}(\theta_G) = -\mathbb{E}_{p(\mathbf{z})p(\mathbf{c})} \left[ \log q_E(\mathbf{z}|\tilde{G}_\tau(\mathbf{z}, \mathbf{c})) \right]$$

# Experiments - ACT

Prompt (Friend)	Response (Friend)	Response (Enemy)
I missed you buddy	Your pic is so cool!	Who doesn't do that
Let's hang out together	Love it been so long	I am horrified by you
Take care I love you	I'm going to miss this	I cannot.

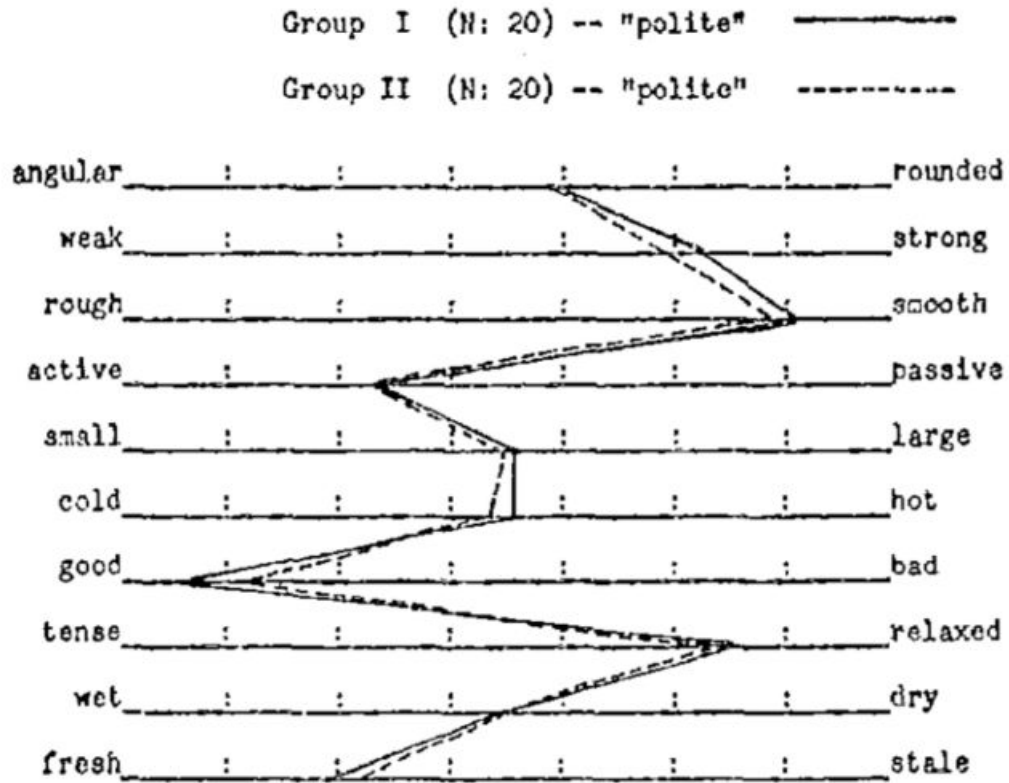
- Constructing sentence given an EPA can be hard (imagine trying to come up with a sentence suitable for (-1, 2, -3).
- Not easy to disentangle emotions from text

# Conclusion

Discussed:

- Continuous vs Discrete Emotion Theory
- Modeling continuous emotions in Seq2Seq framework
- Affective word embeddings, loss functions, diverse decoding
- (Attempting to) Disentangle emotions in latent space
- Use socio-mathematical emotion theory to generate emotional responses

# Osgood's Semantic Differential



## Observations:

- People within one culture answer more or less similarly.
- On average, 50% of variation in semantic differential ratings can be explained by three principal components:

E: good, nice.....bad, awful

P: strong, powerful.....weak, powerless

A: active, excited.....passive, calm