



HANGZHOU
DIANZI
UNIVERSITY



UNIVERSITY
OF
YAMANASHI

AffectI: A Game for Diverse, Reliable, and Efficient Affective Image Annotation

Xingkun Zuo^{1,2}, Jiyi Li², Qili Zhou¹, Jianjun Li¹, Xiaoyang Mao^{1,2}

1. Hangzhou Dianzi University
2. University of Yamanashi

ACM multimedia 2020



- Background
- Related Work
- Proposed Method
- Experiment and Evaluation
- Conclusion and Future work



- **Background**
- Related Work
- Proposed Method
- Experiment and Evaluation
- Conclusion and Future work

Traditional Annotation

Wedding



Cat



Affective Annotation

Happiness
Beautiful

...

Sad
Cute

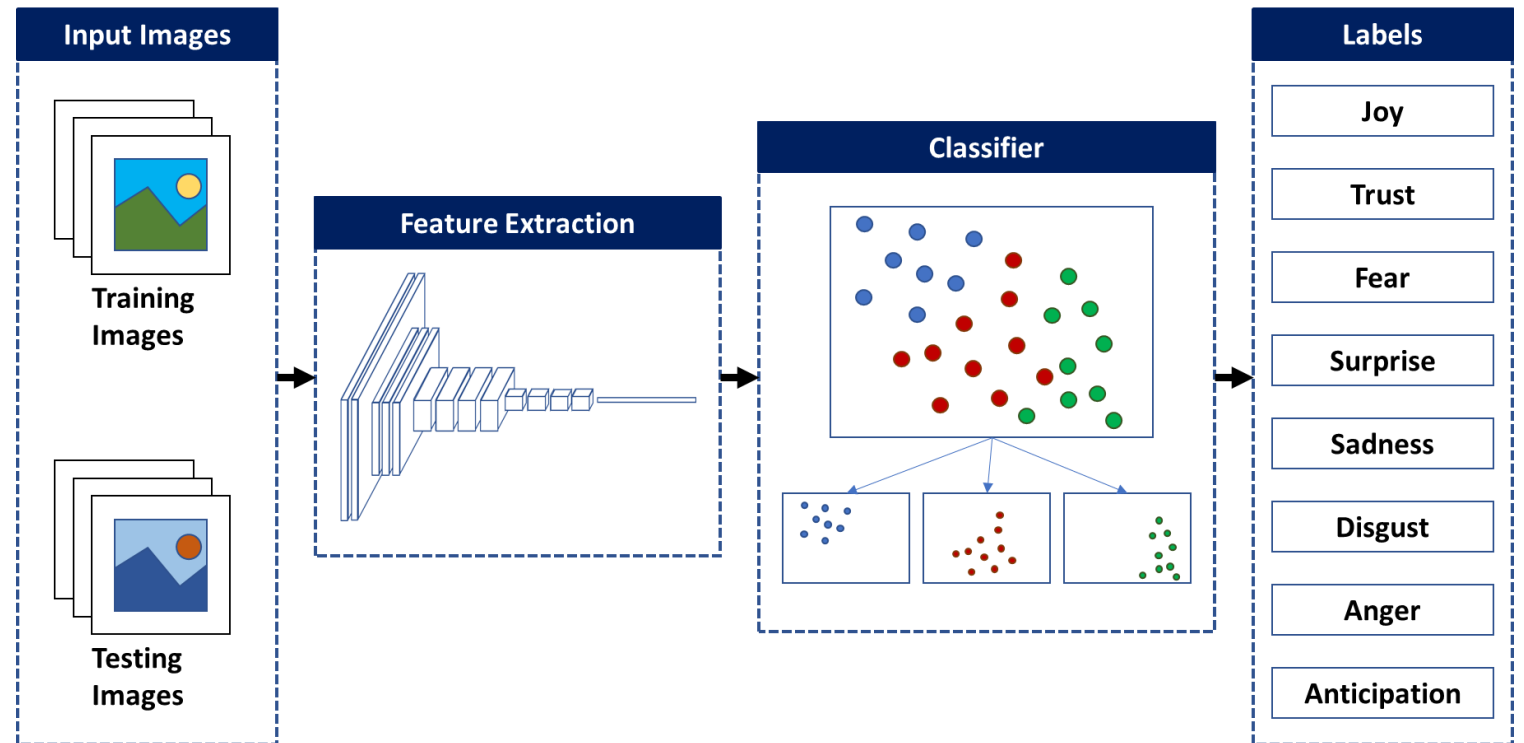
...

Image Retrieval





Affective Image Content Analysis

Annotation is required for collecting training data in machine learning.



AffectI


An online game for labelling images with emotional words

Round:  : 20s Score:  **Start** female

 Your name: **reyaki**


If you don't know how to play, please click tutorial button.






 **Choose player**

Play with system

Your opponent: -

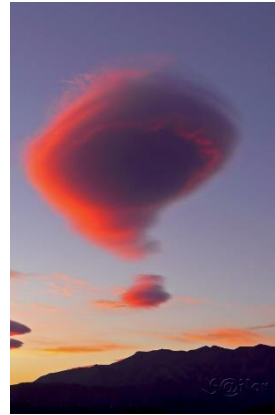
 Number of online players: 1

 Tutorial  Practice  Result Analysis  Emotion Wheel  Chat Box

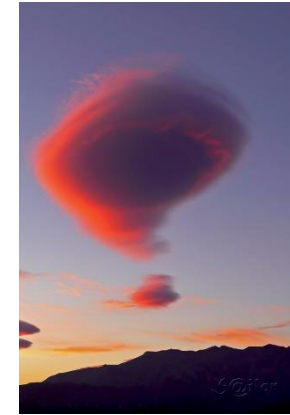


- Background
- **Related Work**
 - Manual labeling
 - Text information
 - Game with A Purpose
- Proposed Method
- Experiment and Evaluation
- Conclusion and Future work

- Manual labeling
 - crowdsourcing
- Text information
- Games with a purpose



Emotion category: **joy**
Sentiment: **positive**



Emotion category: **sadness**
Sentiment: **negative**

Long time manual labeling can easily cause fatigue and affect the reliability of the results

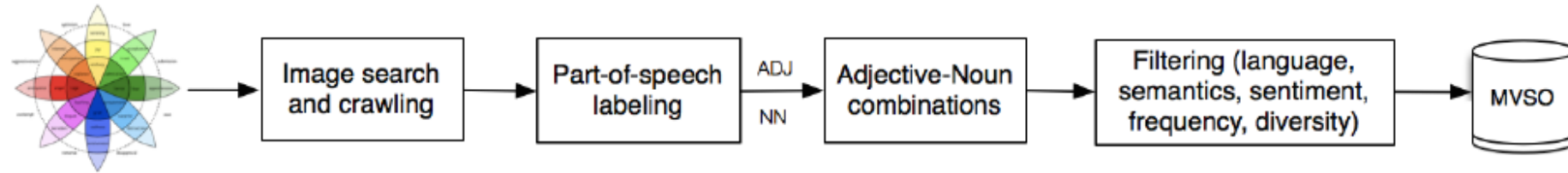
- Manual labeling
- Text information
- Games with a purpose



Description: Yesterday I thought I would spend my short trip very **happily**.

Emotion category: **joy?**

Sentiment: **positive?**



MVSO ——A Large-Scale Multilingual Visual Sentiment Ontology

The quality of the data is usually low.

Game with A Purpose is a human-based computation technique aimed at collecting data as a side effect of game play.

THE ESP GAME

PLAYER 1



GUESSING: **CAR**
GUESSING: HAT
GUESSING: KID
SUCCESS!
YOU AGREE ON CAR

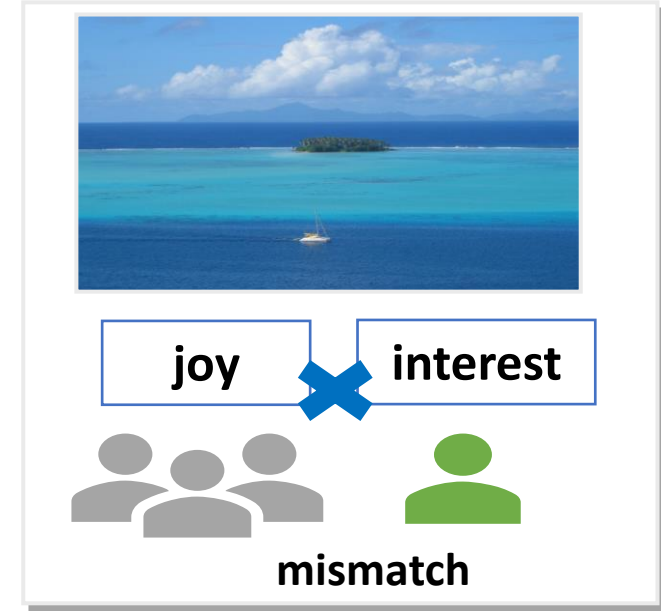
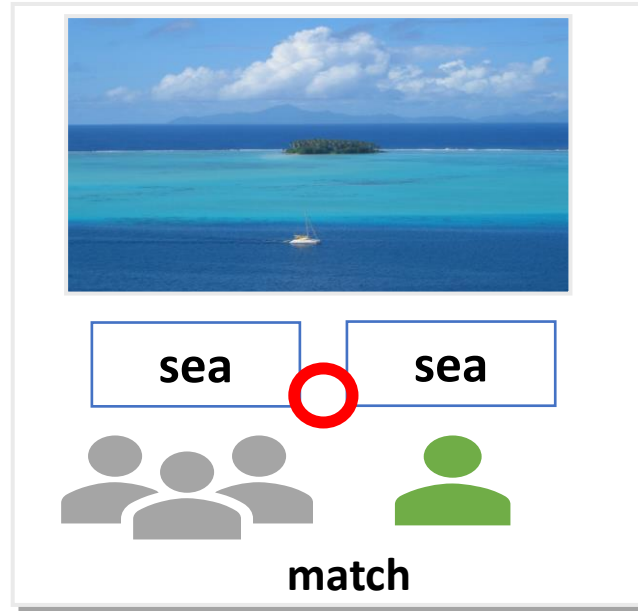
PLAYER 2



GUESSING: BOY
GUESSING: **CAR**
SUCCESS!
YOU AGREE ON CAR

- Manual labeling
- Text information
- Games with a purpose
 - ESP Game
 - KKB Game (Ho et al. 2009)
 - Karido (Bartholomus et al. 2011)

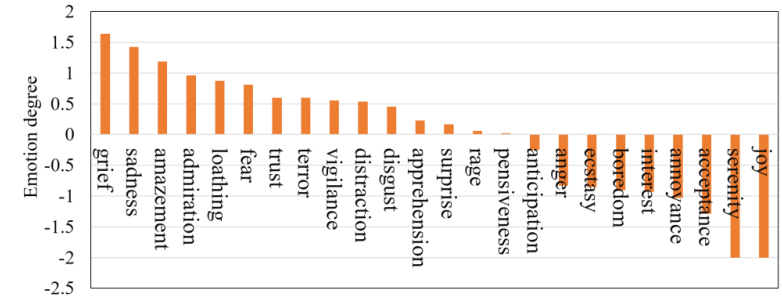
- Manual labeling
- Text information
- Games with a purpose



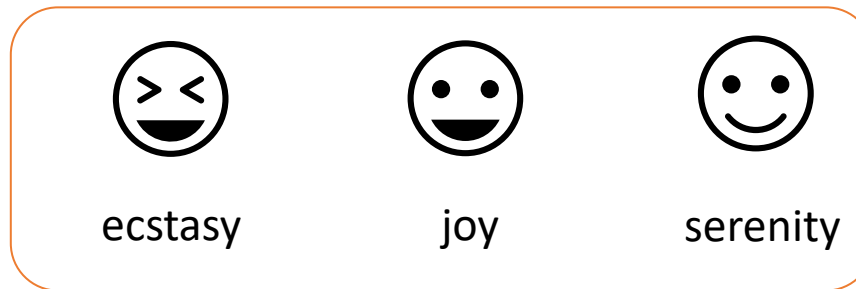
It's difficult to come up with the same emotional words by typing.

Requirement

- Multiple labels
- Emotion distribution



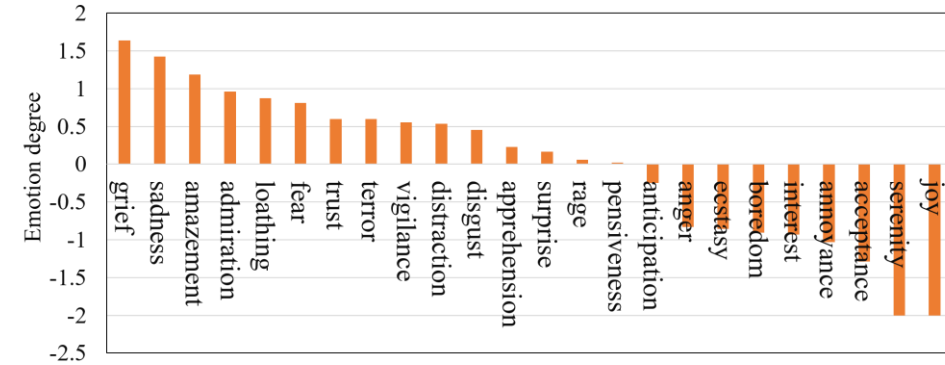
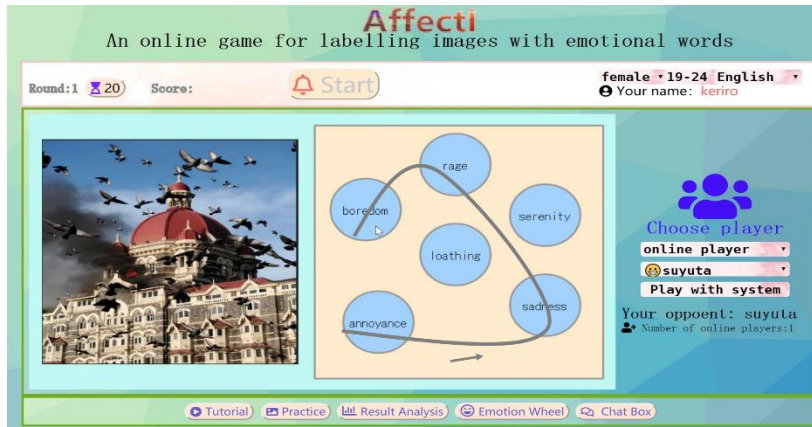
Emotion distribution



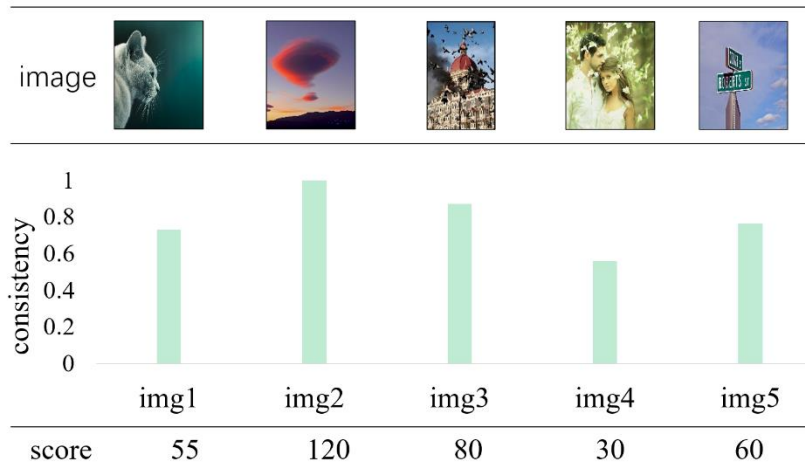
Different degree

Different category

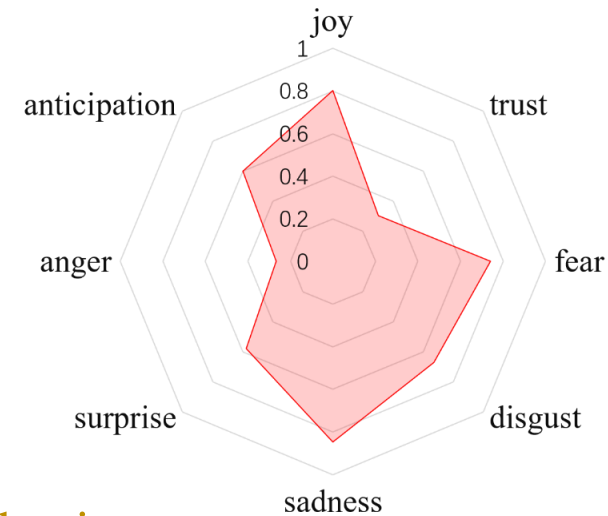
Proposed Method——Three Key Ideas



Selection mechanism



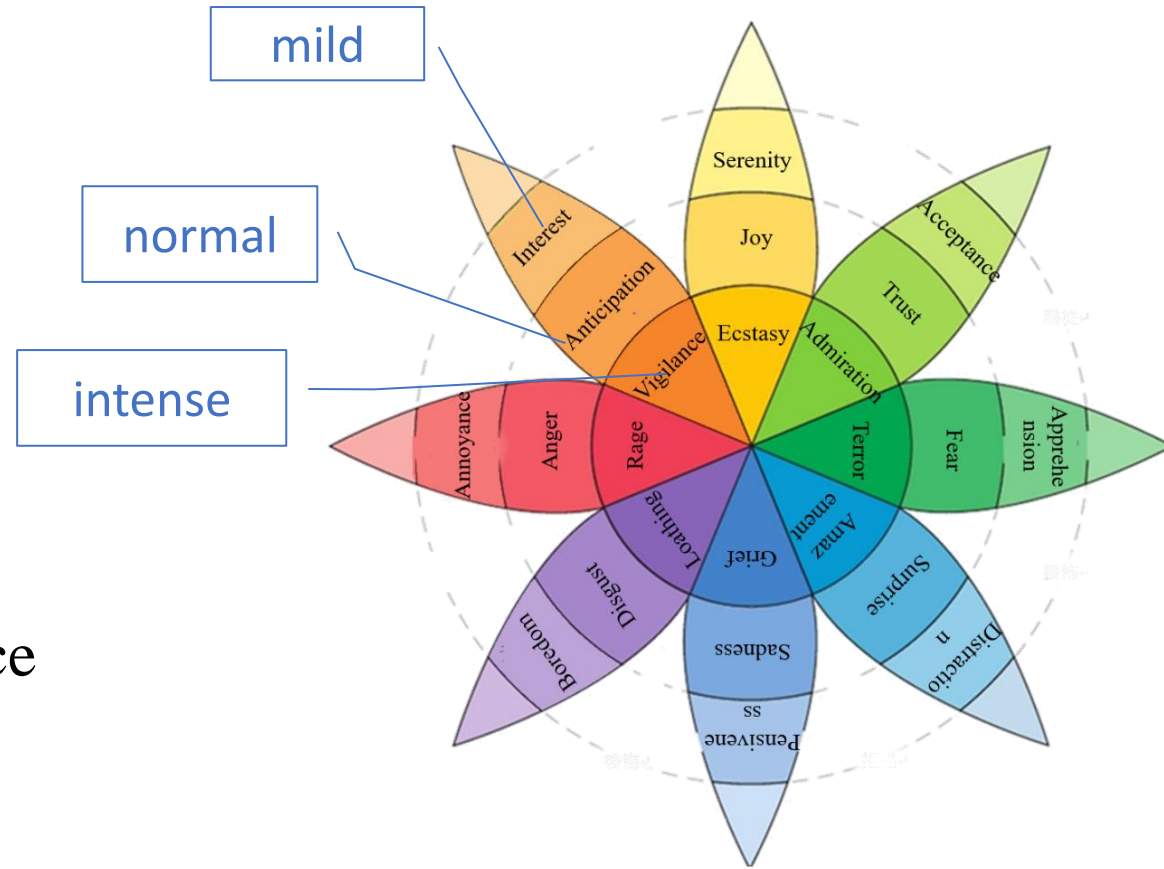
Estimation mechanism



Incentive mechanism

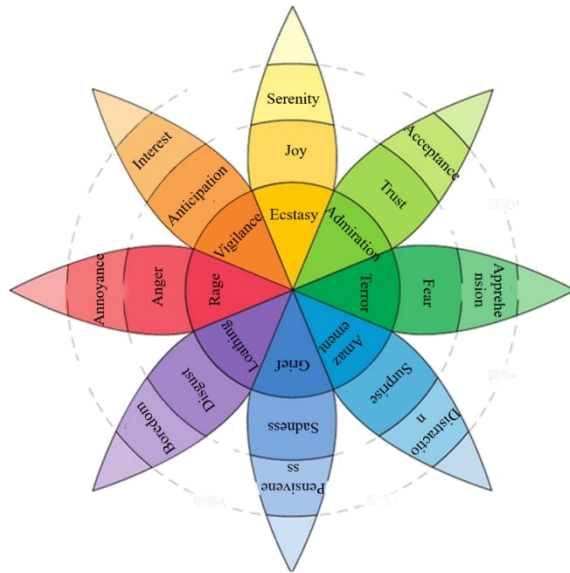
8 categories × 3 scales

intense	normal	mild
ecstasy	joy	serenity
admiration	trust	acceptance
	.	
	.	
	.	



Plutchik's Wheel of emotions

The total number of pairwise comparisons of 24 words is quite large
How to show the words?



Miller's law



The Magical Number 7 + - 2

The number of objects an average human can hold in **short-term memory** is **7 + - 2**.



Our selection mechanism is a hybrid strategy that makes a trade-off between exploitation and exploration.

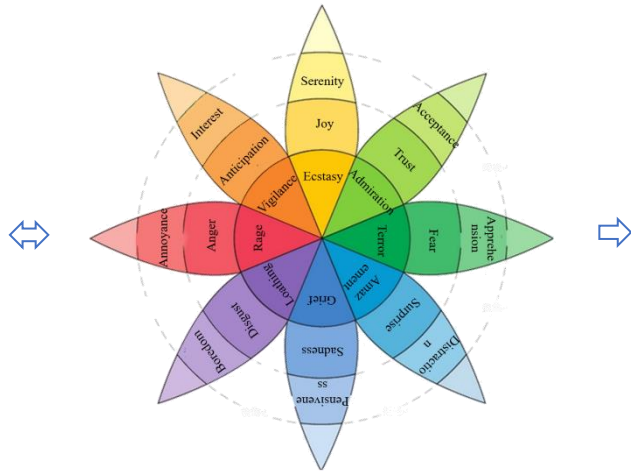
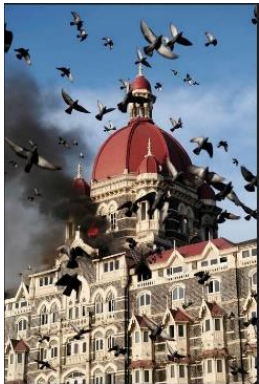
- **Exploitation**

Reconfirm the emotional words that have been labeled to the image more times to obtain a more accurate evaluation.

- **Exploration**

Explore other emotional words that are labeled to the image fewer times. Maybe there have accurate evaluation words in other emotional words.

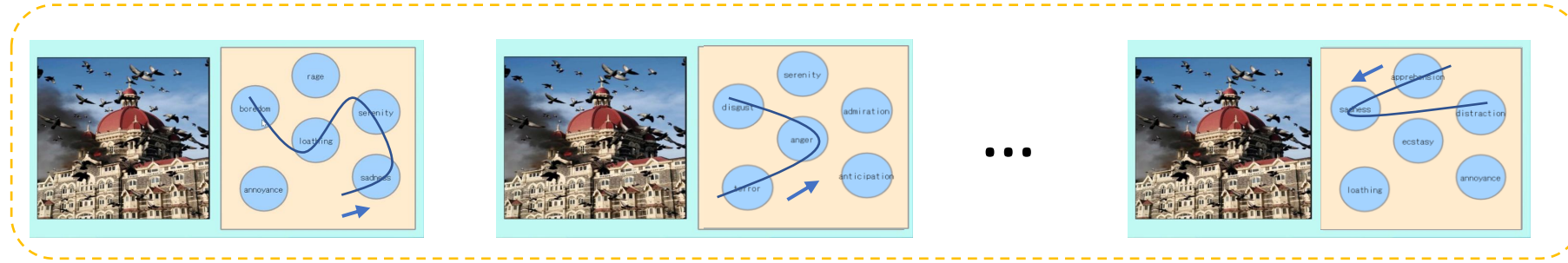
Ensures all emotion words being fairly evaluated.



selected times	displayed times	ρ
Se_1	d_1	$\frac{Se_1}{d_1}$
...
Se_i	d_i	$\frac{Se_i}{d_i}$
...
Se_{24}	d_{24}	$\frac{Se_{24}}{d_{24}}$



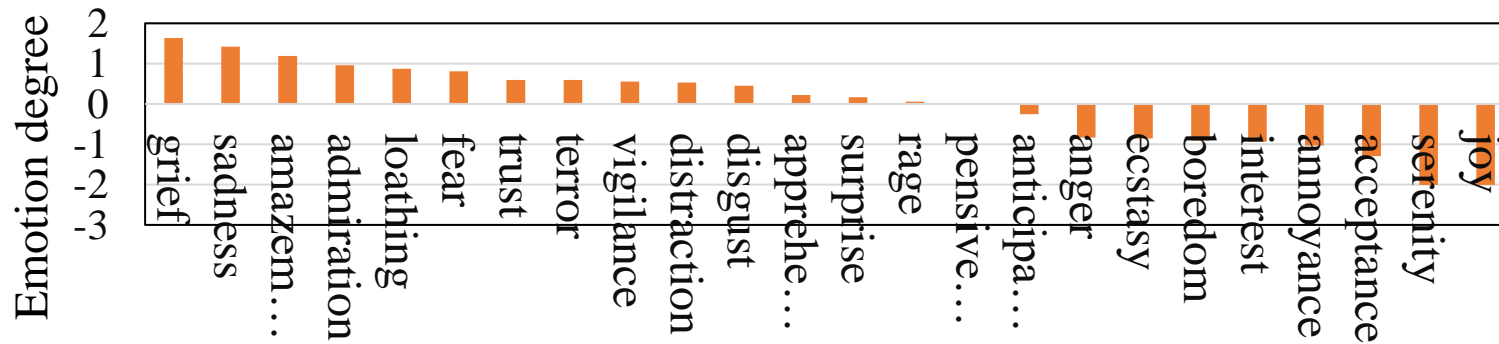
Two: $p > 0.65$
Four: $\min \rho$



A probability model

Relative results

Absolute results



Bradley–Terry model

Estimate the emotion degree s of all emotional words from partial pairwise comparison labels, by minimizing the following objective function:

$$\mathcal{L} = - \sum_{ij} \log(q_{ij}p_{ij} + (1 - q_{ij})(1 - p_{ij}))$$

s_i : the emotion degree of word w_i .

q_{ij} : the observed probability that word w_i precedes word w_j

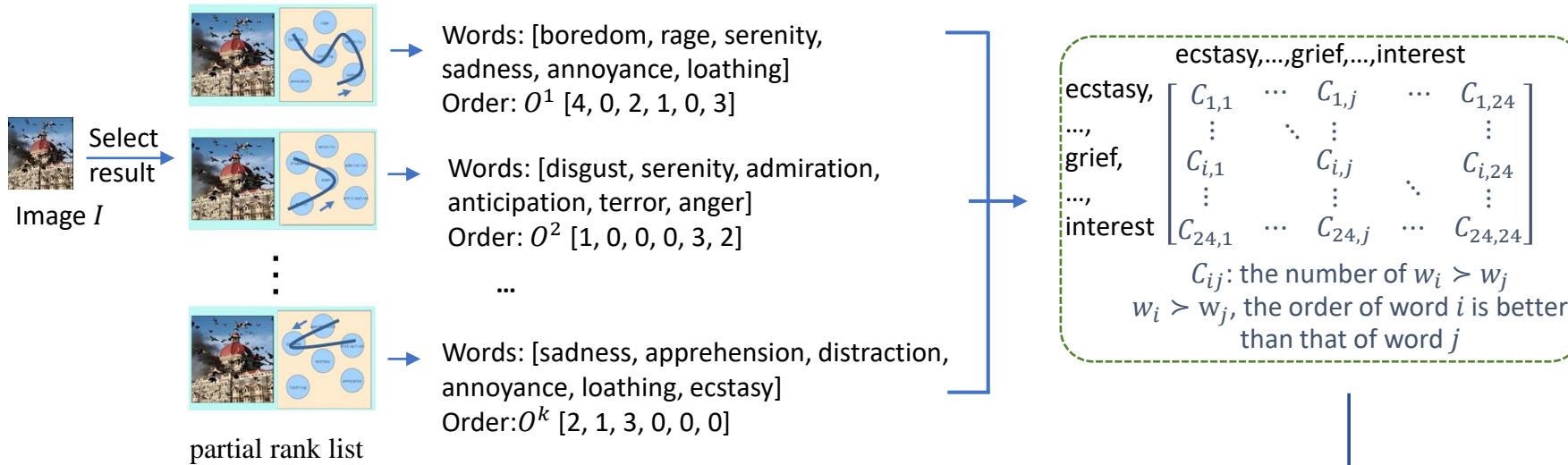
p_{ij} : the predicted possibility that word w_i precedes word w_j

$$p_{ij} = \frac{e^{s_i}}{e^{s_i} + e^{s_j}} = \frac{1}{1 + e^{-(s_i - s_j)}}$$

The objective function can be rewritten as:

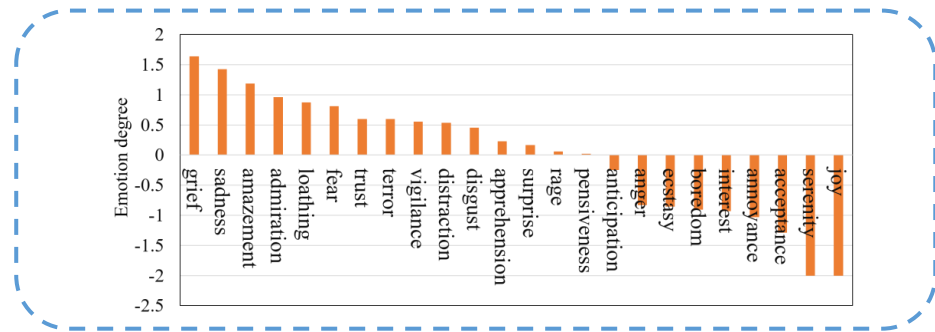
$$\mathcal{L} = - \sum_{ij} \log\left(q_{ij} \frac{1}{1 + e^{-(s_i - s_j)}} + (1 - q_{ij}) \frac{1}{1 + e^{-(s_j - s_i)}}\right)$$

Estimate the Emotion Degree



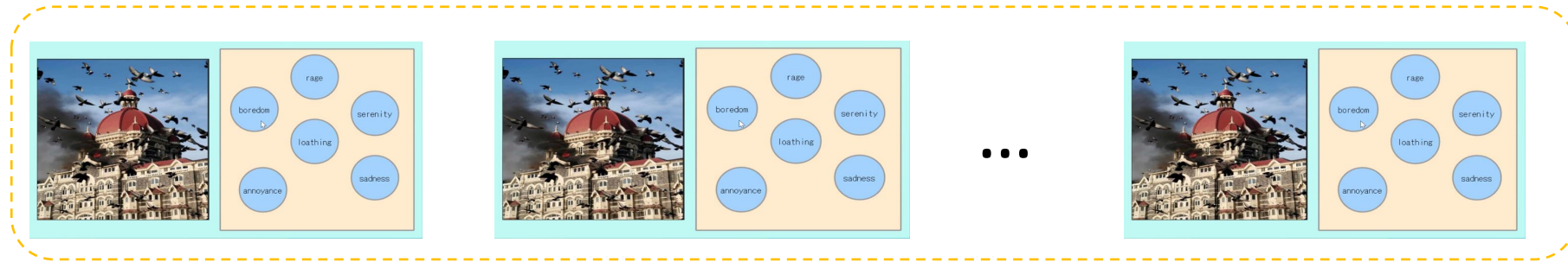
$$p(w_i > w_j) = \frac{e^{s_i}}{e^{s_i} + e^{s_j}} = \frac{1}{1 + e^{-(s_i - s_j)}}$$

s_i : the emotion degree of word w_i .



What is the Matrix C?

All partial rank lists will be transformed into pairwise preference comparison.

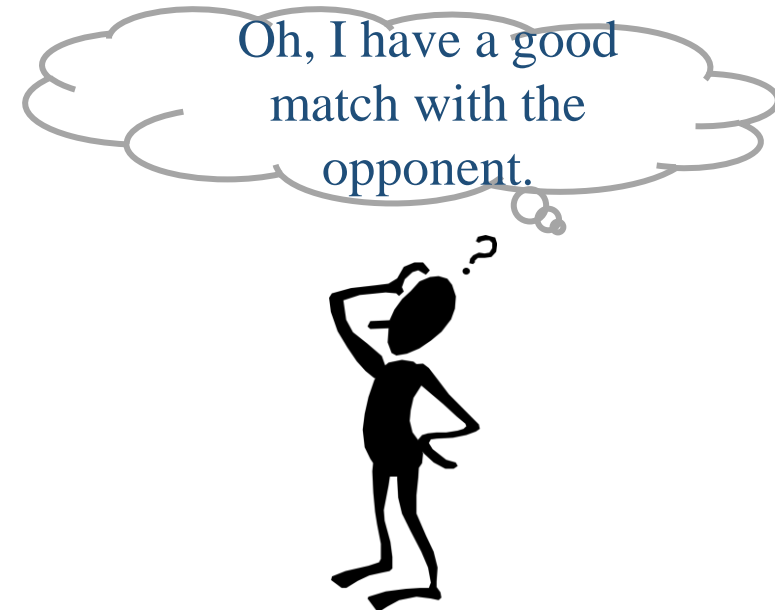
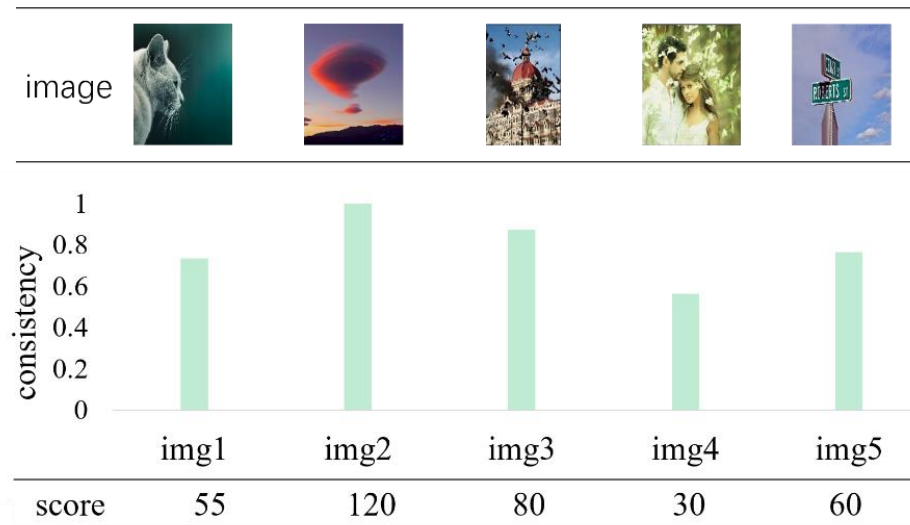


The pairwise preference comparison $w_i \succ w_j$ means the order of the word i is better than the word j .

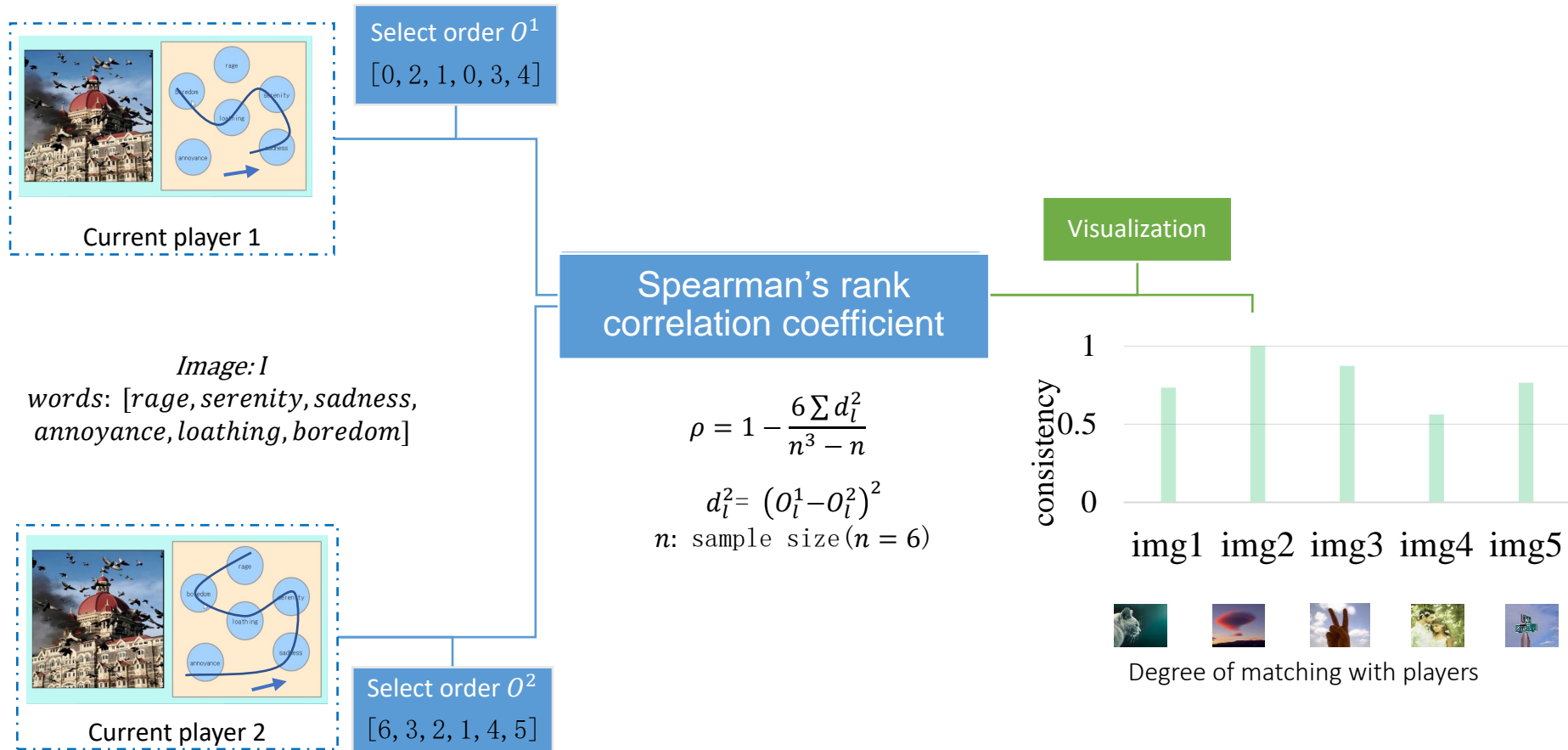
For each subset, if the word i is labeled for the image, but the word j is not labeled, it also means the order of word i is better than word j .

Comparison between the current player and her opponent

To encourage the player to provide labels with high quality.

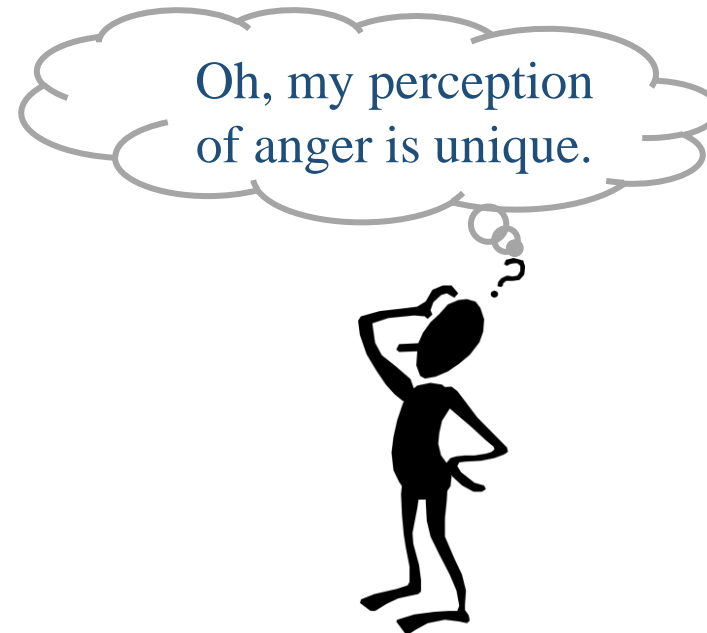
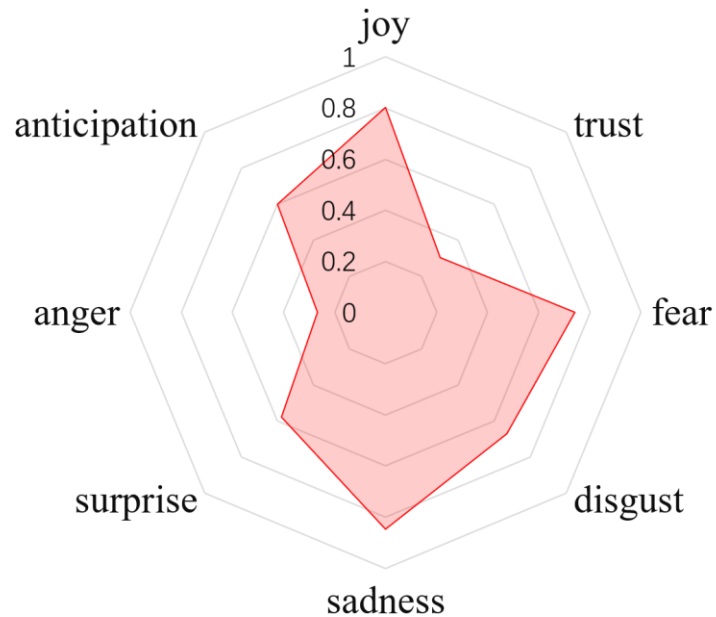


Opponent-based Incentive (OI)

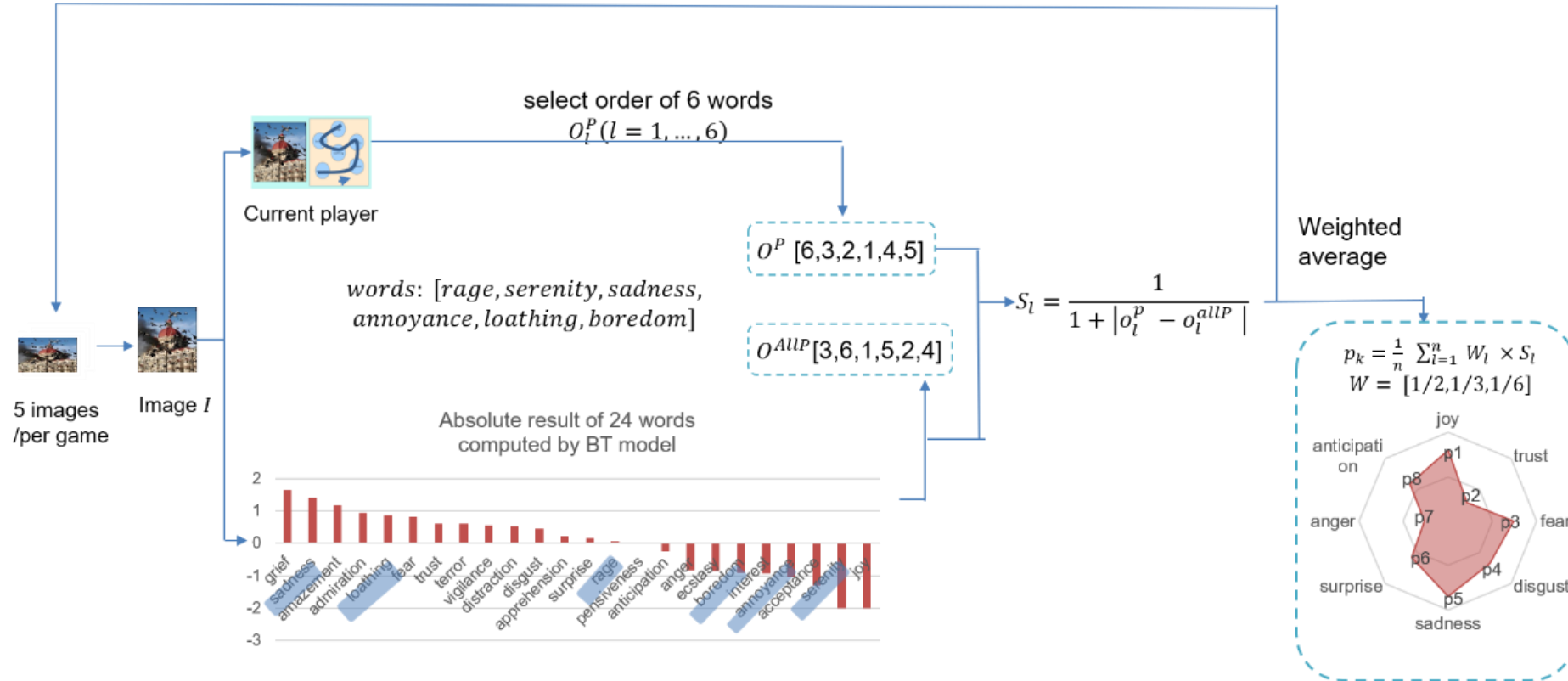


Comparison between the current player and all past players

To encourage more players to choose the labels representing their own subjective emotional perception.



Past Players-Based Incentive (PPI)



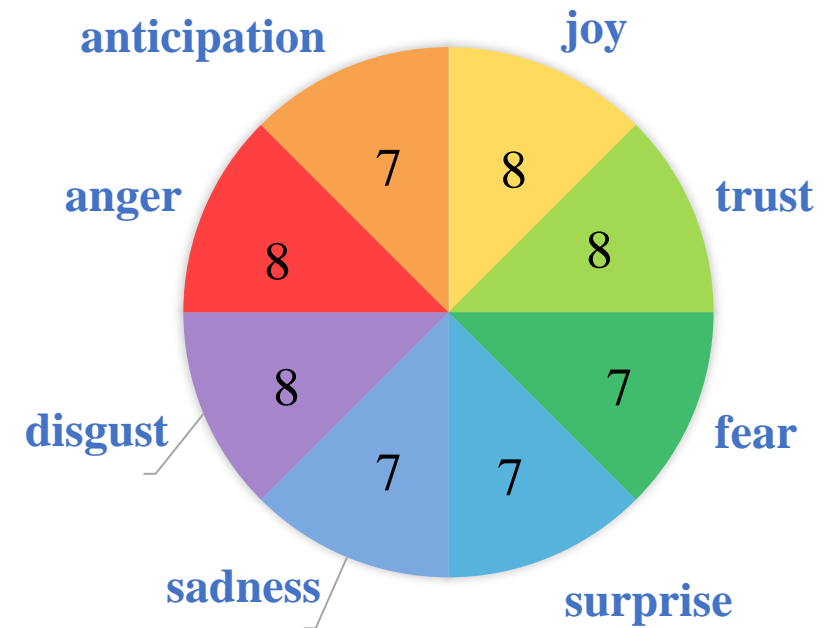


- Background
- Related Work
- Proposed Method
- **Experiment and Evaluation**
- Conclusion and Future work

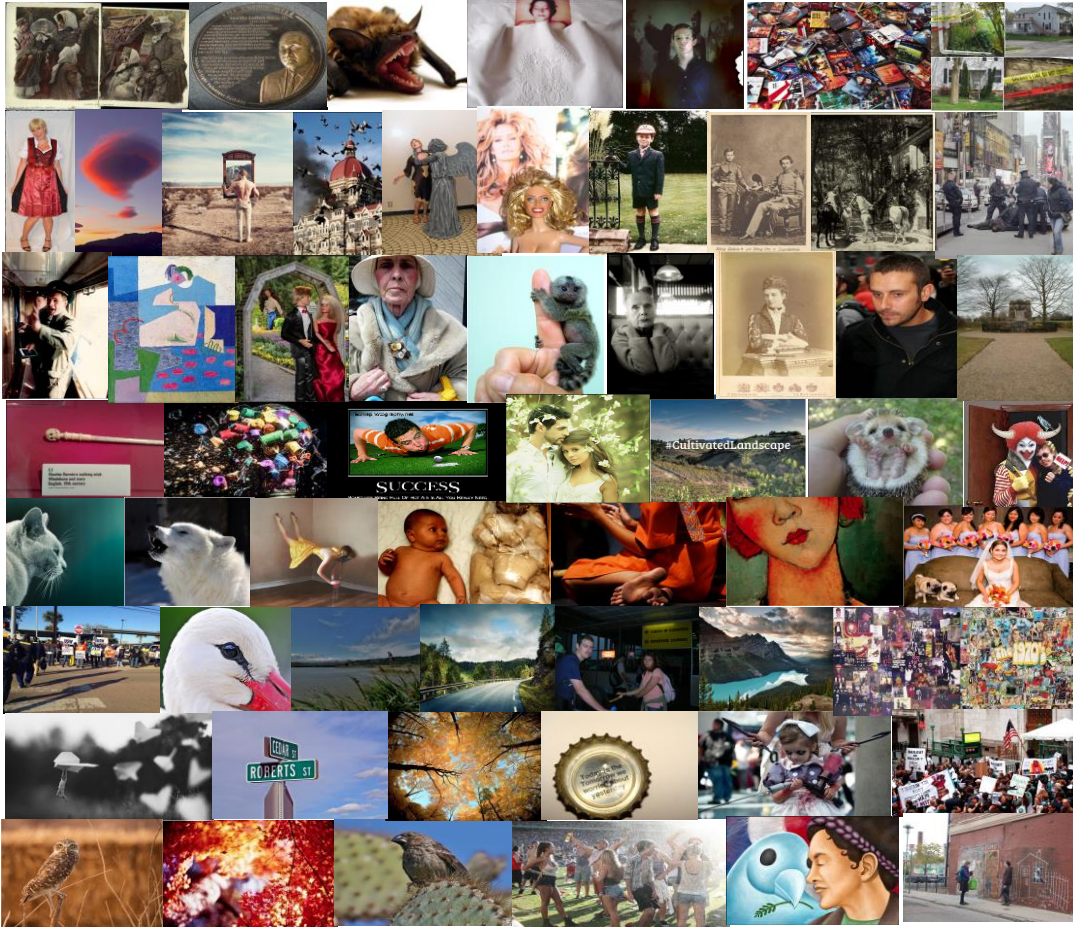
60 images from MVSO dataset (sub_MVSO)

Each category: 7~8 highly viewed images on flicker

Each strength level: 1~3 images



Experiment—Results



- “Our-PPI”: without PPI
- “Our+PPI”: with PPI

Number of images

60

Participants

Age: 22-35

Number of
unique player
IDs

Our-PPI: 163

Our+PPI: 67

Labeled times

Our-PPI: 1,892

Our+PPI: 710



Evaluation——What to Evaluate



Evaluation——Diversity

Can the system successfully collect diverse emotional labels with different emotion degrees compared to the existing work?



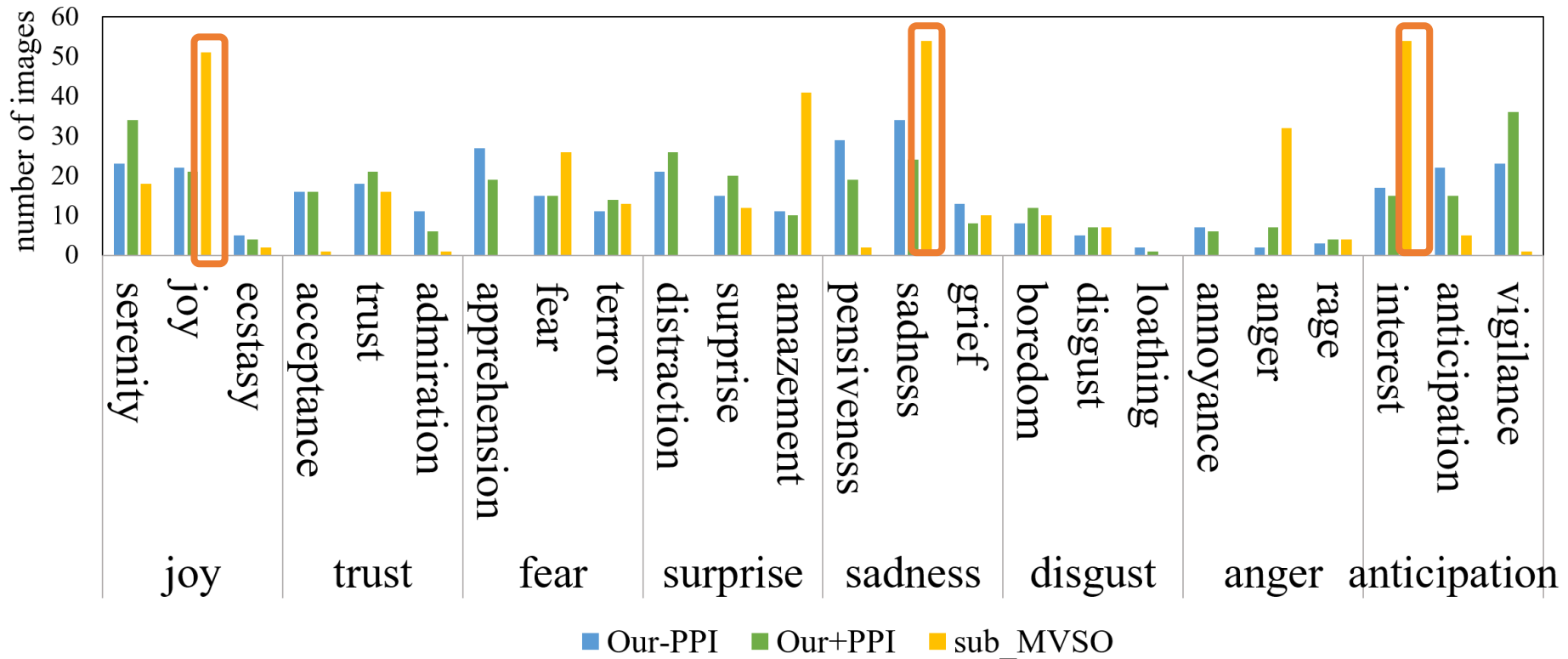
Fear, anger,
interest,



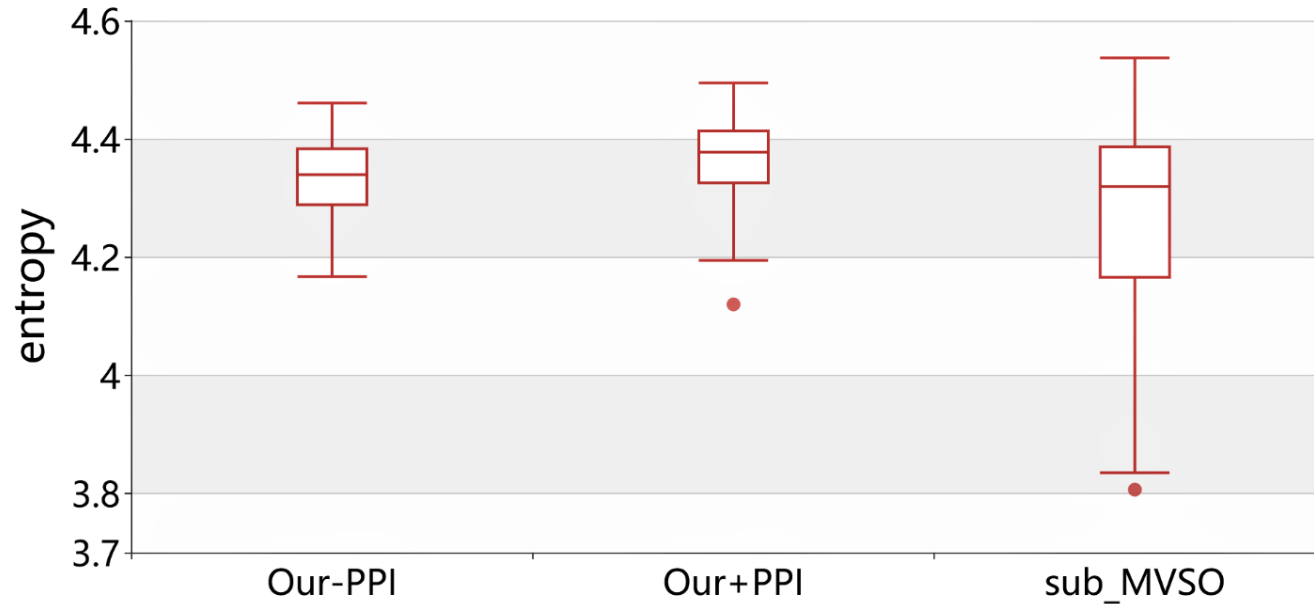
Distraction, joy,
serenity,



Comparing the frequency of the top-6 emotion words annotated for 60 images.



An entropy-based measurement.



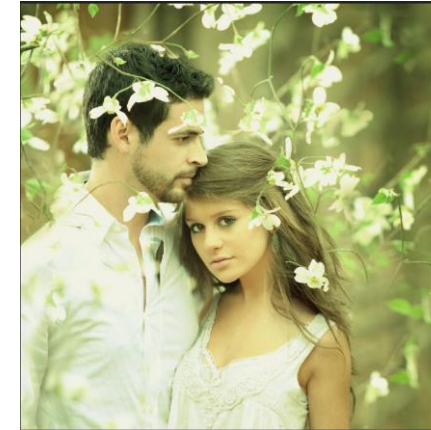
The average entropy of the distribution of the 24 emotions for 60 images.

Evaluation——Reliability

What is the quality of the diverse labels collected by our system?



sadness
amazement
joy
interest
serenity
anger



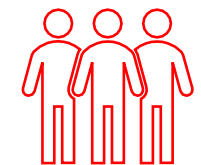
joy
anticipation
interest
admiration
serenity
trust

Existing work



agree

Our work



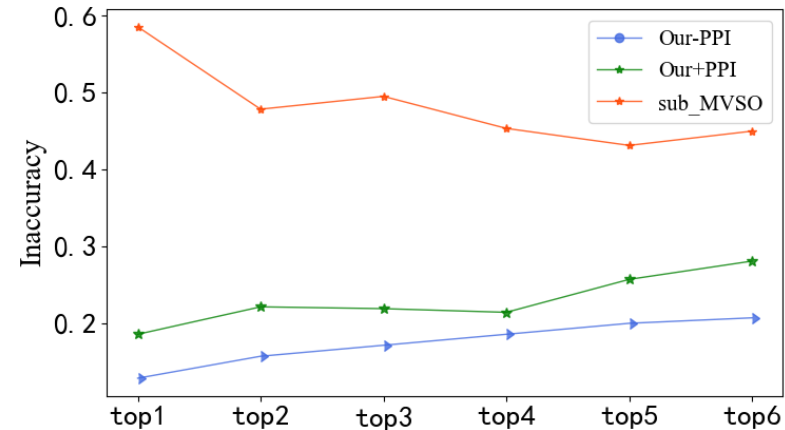
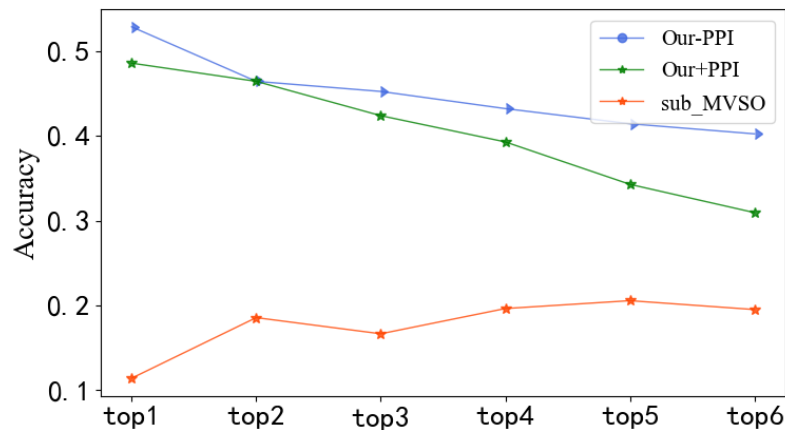
agree

Images: 10,

Evaluator: 7 participants /per group, 3~4 images/per group,

Words: top-6 emotion words obtained by all system,

Judgments: inaccurate, neutral, and accurate



Comparison of the average votes of accurate and inaccurate for the top- k ($k \leq 6$) emotional labels.

What is the player experience when using AffectI?



Manual labeling experiment involving 28 participants for comparison.

Images: 15 images / per group

Participants: 7 / per group

Q: Please select the emotional word that you think the image best represents or the emotional word caused by the image.



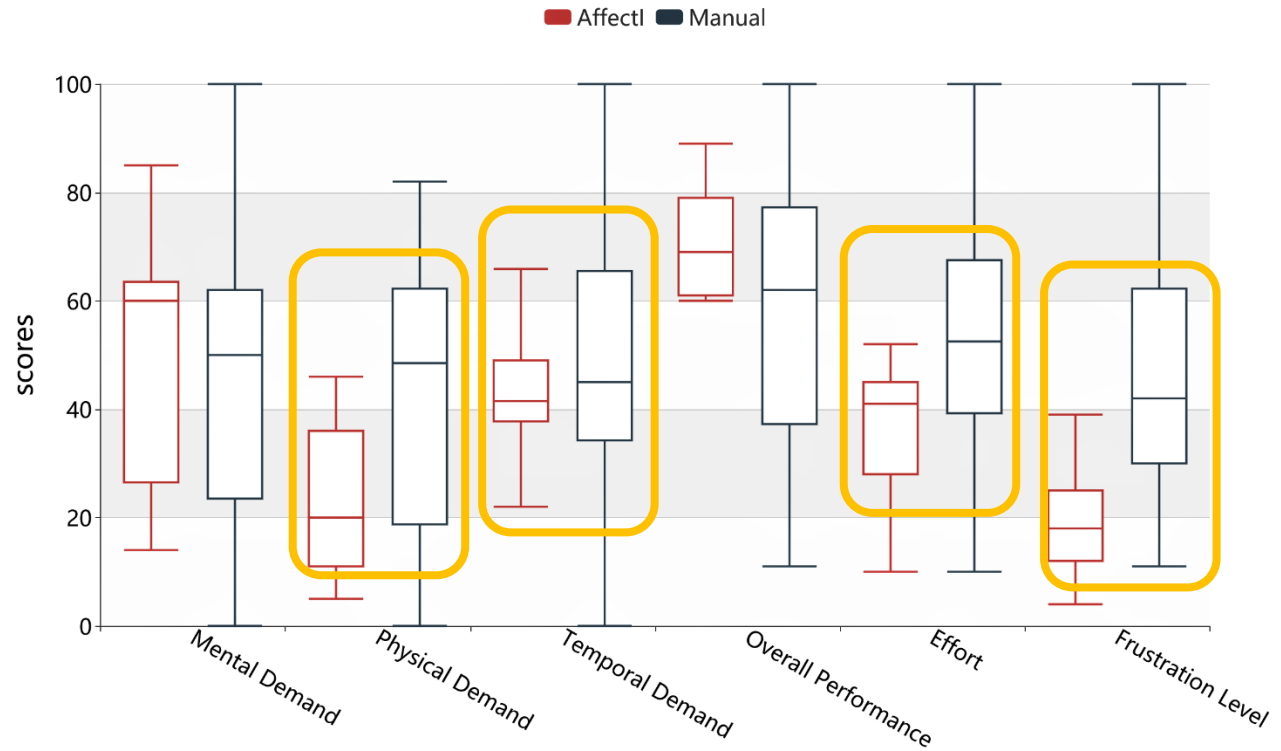
- | | | |
|-------------------------------------|---------------------------------------|---------------------------------------|
| <input type="checkbox"/> Ecstasy | <input type="checkbox"/> Joy | <input type="checkbox"/> Serenity |
| <input type="checkbox"/> Admiration | <input type="checkbox"/> Trust | <input type="checkbox"/> Acceptance |
| <input type="checkbox"/> Terror | <input type="checkbox"/> Fear | <input type="checkbox"/> Apprehension |
| <input type="checkbox"/> Amazement | <input type="checkbox"/> Surprise | <input type="checkbox"/> Distraction |
| <input type="checkbox"/> Grief | <input type="checkbox"/> Sadness | <input type="checkbox"/> Pensiveness |
| <input type="checkbox"/> Loathing | <input type="checkbox"/> Disgust | <input type="checkbox"/> Boredom |
| <input type="checkbox"/> Rage | <input type="checkbox"/> Anger | <input type="checkbox"/> Annoyance |
| <input type="checkbox"/> Vigilance | <input type="checkbox"/> Anticipation | <input type="checkbox"/> Interest |

Comparison with the results of 15 randomly selected AffectI players

NASA Task Load Index

Hart and Staveland's NASA Task Load Index (TLX) method assesses work load on five 7-point scales. Increments of high, medium and low load on five 7-point scales. Increments of high, medium and low estimates for each point result in 21 gradations on the scales.

Name	Task	Date
Mental Demand How mentally demanding was the task?		
Physical Demand How physically demanding was the task?		
Temporal Demand How hurried or rushed was the pace of the task?		
Performance How successful were you in accomplishing what you were asked to do?		
Effort How hard did you have to work to accomplish your level of performance?		
Frustration How insecure, discouraged, irritated, stressed, and annoyed were you?		



NASA-TLX evaluation results



- Background
- Related Work
- Proposed Method
- Experiment and Evaluation
- **Conclusion and Future work**



Conclusion

We proposed a novel affective image annotation system, AffectI, which can efficiently collect high-quality and diverse emotional image dataset.

Future Work

- To enlarge the dataset
- To estimate the personalized emotion degrees of players
- To analyze how emotion perception on images is related to other factors such as age, gender, and language.



HANGZHOU
DIANZI
UNIVERSITY



UNIVERSITY
OF
YAMANASHI

Thanks for your attention!