



UMASS  
AMHERST



# Reasoning about Fine-grained Attribute Phrases using Reference Games

Jong-Chyi Su\* Chenyun Wu\* Huaizu Jiang Subhransu Maji  
University of Massachusetts, Amherst

## Motivation

Different from fixed attributes designed by experts, we collect **attribute phrases** describing differences between two images. They can be better generalized to new domains, and handle open-ended descriptions from non-expert users.

**Reference Game:**  
We learn to **describe** and **ground** attribute phrases to images using reference game between a **speaker** and a **listener**.

### Dataset:

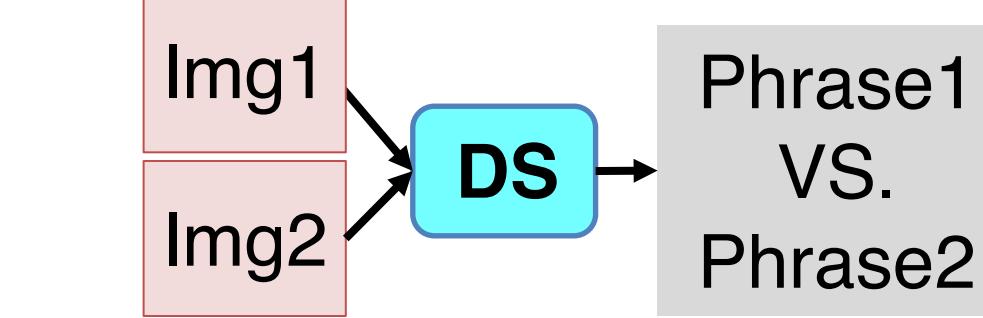
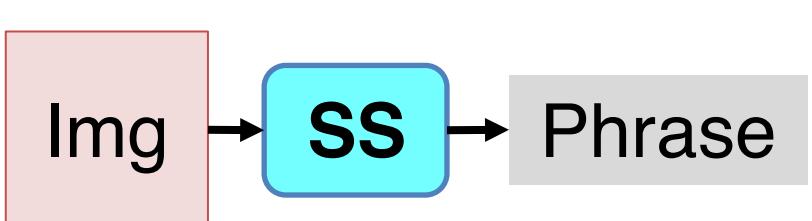
We collect a dataset based on OID<sup>[1]</sup> dataset, containing 9400 image pairs and 5 descriptions per pair.



## Approach

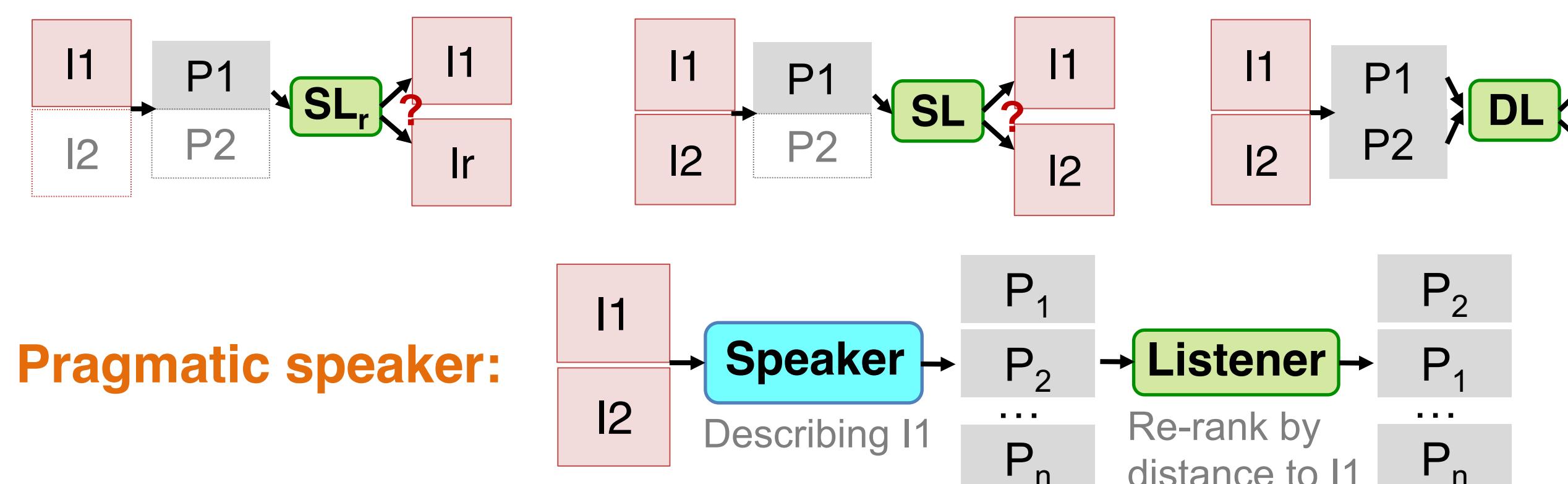
### Speaker [SS, DS]:

We follow the Show-and-tell image captioning model<sup>[2]</sup>.



### Listener [SL<sub>r</sub>, SL, DL]:

We train model to measure the similarity between the phrases and images in a common embedded space.



## Quantitative Result

### →: Listener

The accuracy of picking the target image out of a pair

	Accuracy (%)				
	SL <sub>r</sub>		SL		Human
Top	Test*	Test	Test*	Test	Test*
SS	1	84.0	79.8	83.0	81.7
	5	80.0	79.2	78.0	80.6
	10	78.0	78.9	76.6	80.0
DS	1	94.0	92.8	92.0	82.0 (88.5)
	5	91.2	90.3	91.2	80.2 (86.7)
	10	88.6	88.8	90.0	90.5

Input	Speaker	Listener	Val	Test
P <sub>1</sub>	Human	SL <sub>r</sub>	82.7	84.2
		SL	85.3	86.3
P <sub>1</sub> vs. P <sub>2</sub>	Human	DL	88.7	88.9
		2×SL	89.6	89.3

Input	Speaker	Human listener accuracy (%)	
		Reranker listener	SL <sub>r</sub>
SS	Top	None	SL <sub>r</sub>
	1	68.0 (77.0)	94.0 (96.0)
	5	64.2 (74.1)	82.6 (88.3)
DS	7	63.1 (72.8)	74.3 (82.0)
	1	82.0 (88.5)	95.0 (96.5)
	5	80.2 (86.7)	90.0 (93.3)
7	79.1 (85.6)	86.7 (91.5)	86.1 (91.1)

### ↑: Speaker

Use speaker models to generate phrases, evaluated by different listeners

### ↑: Pragmatic speaker

The accuracy from human listener of the top-k phrases generated by speakers and re-ranked by listeners

## Example Output

### →: Speaker

DS sees two images  
SS sees only the image in the green box  
**DS is better than SS!**

### ↓: Pragmatic speaker

**Re-ranking improves!**



**Ground Truth:**  
1) small size **VS** large size  
2) single seat **VS** more seated  
3) facing left **VS** facing right  
4) private **VS** commercial  
5) wings at the top **VS** wings at the bottom

**DS:**

- 1) private plane **VS** commercial plane ( $p=0.3338$ )
- 2) private **VS** commercial ( $p=0.1648$ )
- 3) small plane **VS** large plane ( $p=0.0701$ )
- 4) facing left **VS** facing right ( $p=0.0355$ )
- 5) short **VS** long ( $p=0.0250$ )
- 6) white **VS** red ( $p=0.0228$ )
- 7) high wing **VS** low wing ( $p=0.0184$ )
- 8) small **VS** large ( $p=0.01775$ )
- 9) glider **VS** jetliner ( $p=0.0170$ )
- 10) white and blue color **VS** white red and blue color ( $p=0.0159$ )

**SS:**

- ✓ passenger plane
- ? white
- ✓ large
- ✓ facing right
- ✓ jet engine
- ✓ commercial plane
- ? facing left
- ? \_UNK
- ? on the ground
- ✓ large size
- ✓ turbofan engine
- ✓ on concrete
- ✓ t tail
- ✓ passenger plane
- ? on the ground
- ✓ jet engine
- ✓ twin engine
- ✓ on concrete
- ✓ multi seater
- ✓ t tail
- ✓ white and red
- ✓ white colour with red stripes

**SS + SL<sub>r</sub>:**

- ✓ commercial plane
- ✓ facing right
- ✓ jet engine
- ✓ large size
- ✓ turbofan engine
- ✓ on concrete
- ✓ t tail
- ✓ passenger plane
- ? on the ground
- ✓ jet engine
- ✓ twin engine
- ✓ on concrete
- ✓ multi seater
- ✓ t tail
- ✓ white and red
- ✓ white colour with red stripes

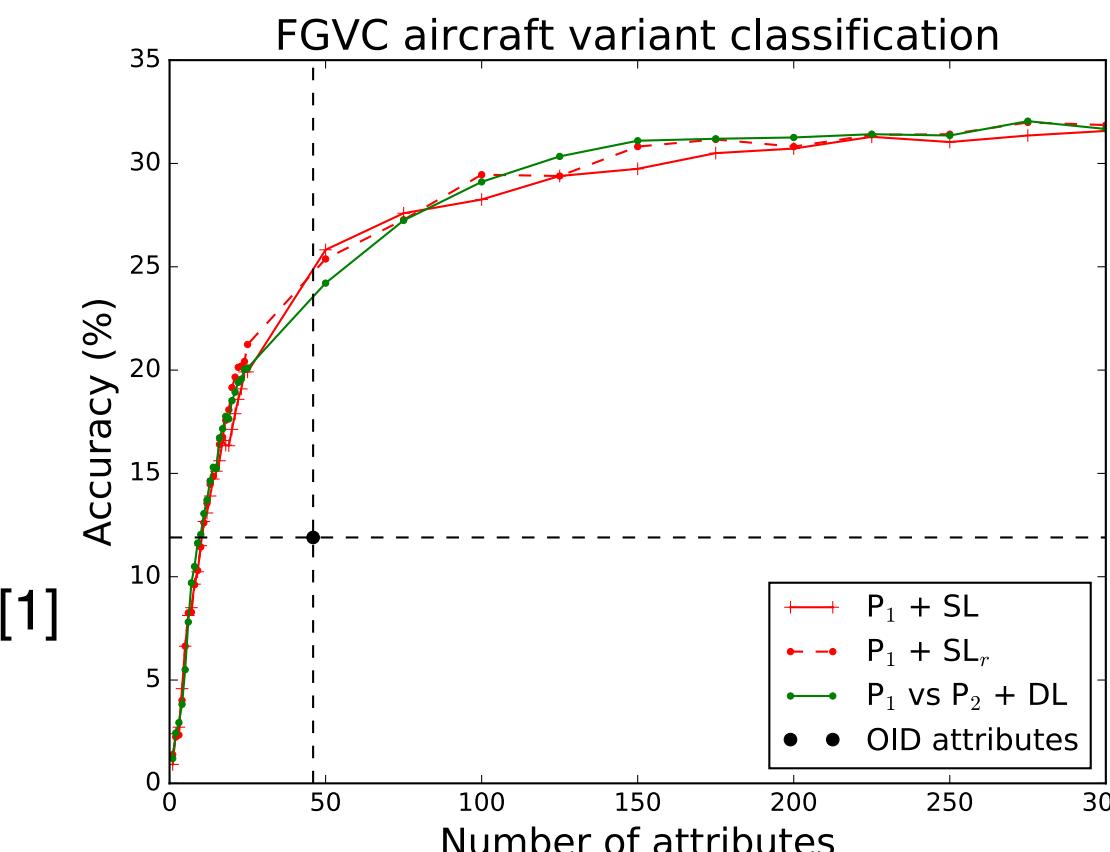
**DS + SL<sub>r</sub>:**

- ✓ commercial plane
- ✓ facing right
- ✓ jet engine
- ✓ large size
- ✓ turbofan engine
- ✓ on concrete
- ✓ t tail
- ✓ passenger plane
- ? on the ground
- ✓ jet engine
- ✓ twin engine
- ✓ on concrete
- ✓ multi seater
- ✓ t tail
- ✓ white and red
- ✓ white colour with red stripes

## Fine-grained Classification

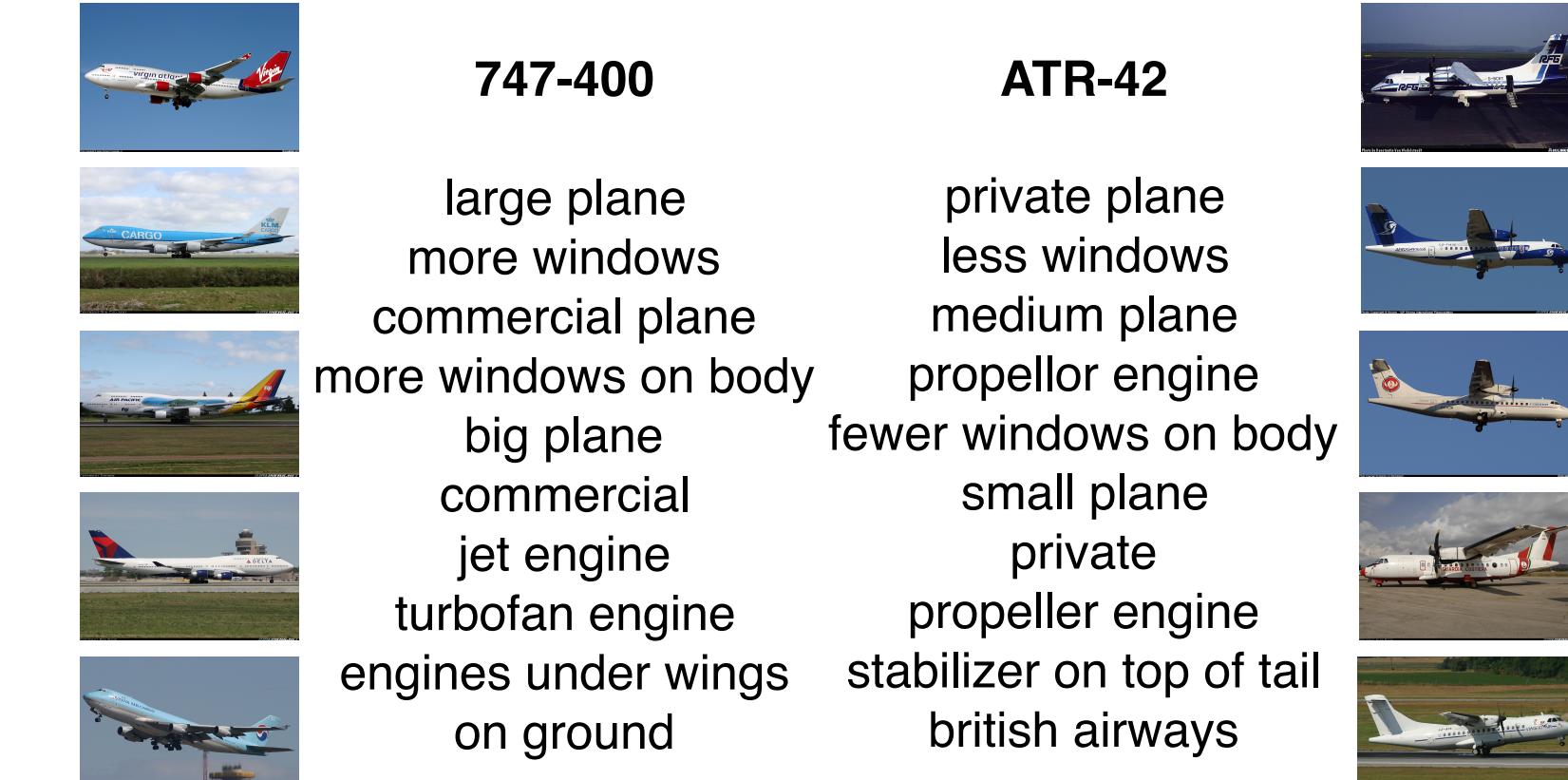
Use attributes as features for fine-grained classification on FGVC-aircraft<sup>[3]</sup> dataset:

1. Use our listener model to get scores between the *top-k most frequent attribute phrases* and the image
  2. Expert-designed 46 attributes from OID<sup>[1]</sup>
- Our attribute phrases outperform by **20%**.

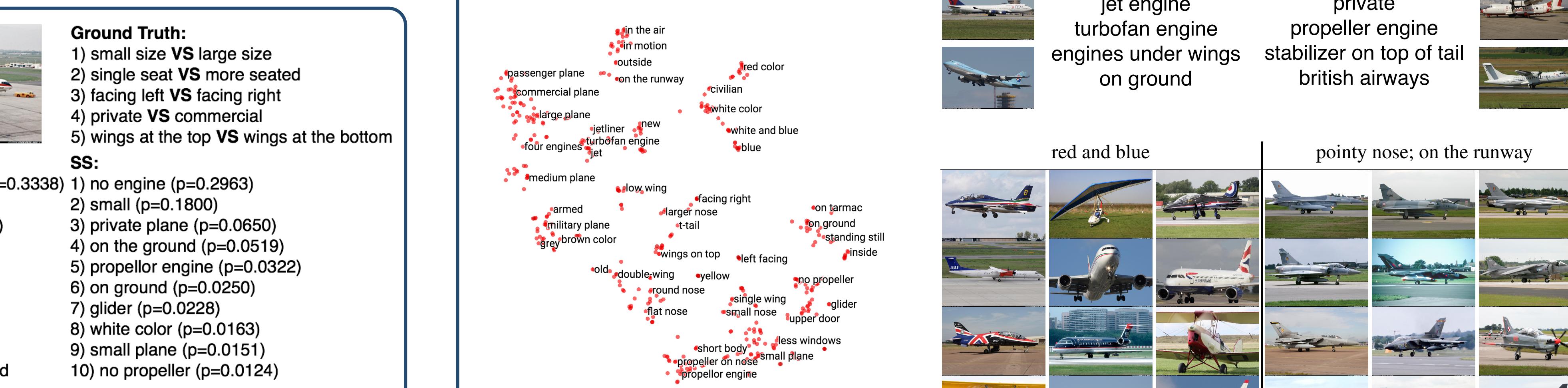


## Visualization & Application

### →: Set-wise attributes



### ↓: t-SNE Embedding of attribute phrases



### →: Image retrieval

Top-18 retrieved images in the test set, ranked by the listener



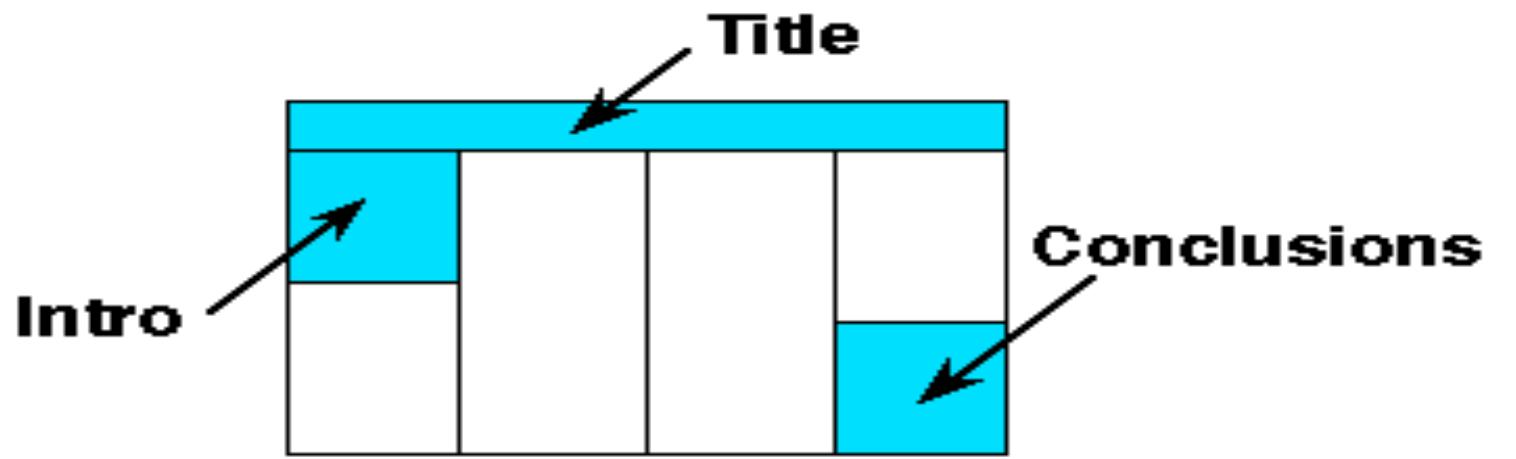
[1] Vedaldi et al., Understanding objects in detail with fine-grained attributes, CVPR, 2014.

[2] Vinyals et al., Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge, TPAMI, 2016.

[3] Maji et al., Fine-grained visual classification of aircraft. arXiv:1306.5151, 2013.

## Section 1 (sizes):

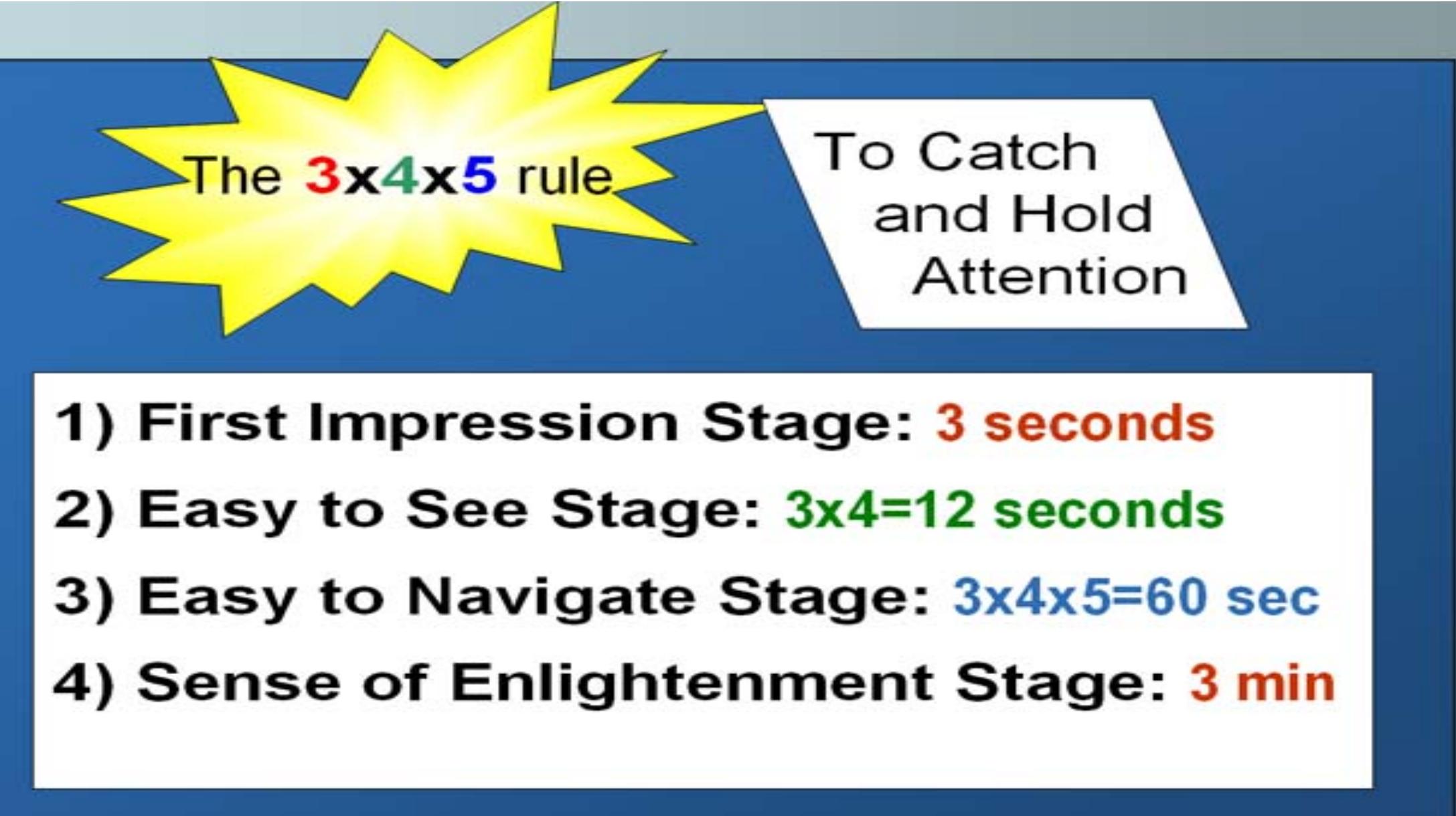
- Posters boards are 48" tall and 96" wide, but we recommend you leave a little border since you may not be able to pin at the vertical edge. Since PowerPoint does not let one define such a large paper size, this template is designed to be printed at 200%, yielding a 46" x94" poster. You can scale it up or down a bit (e.g. 42" is a common paper size at FexEd). Note there is no direct international A0.. A1 equivalent. The poster size is approximately three A0 boards next to each other, i.e., each column in this example is about one A0 board.
- Ideally you want to keep it very readable: this is not your paper, it is a poster. 32pt here (64 final printing) is good for most text:
  - Sub-bullets are 28 here (56 final)
  - Don't use smaller than 24pt in this template (which is 48pt in final printing at 200%)
  - Insert plenty of graphics and any math you need
- When inserting graphics or equations, keep the resolution high (remember this will be printed at 200%). If you can see blocking artifacts at 400% magnification in PowerPoint, consider finding better graphics. This is an example of BAD/LOW RES GRAPHICS



- Leave enough margin for pushpin and remember many big plotters cannot get within .5" of the actual paper edge.
- You are free to use colored backgrounds and such but they generally reduce readability.
- You are free to use what ever fonts you like.
  - San Serif fonts like Arial are more readable from a distance,
  - Serif fonts like times may look more consistent with your mathematics

## Section 2 (layout):

- Remember the poster session will be crowded so design the poster to be read in columns so people can read what is in front of them and move left to right to get the whole story.
- The poster should use photos, figures, and tables to tell the story of the study. For clarity, present the information in a sequence that is easy to follow.
- There is often way too much text in a poster - there definitely is in this template! Posters primarily are visual presentations; the text should support the graphics. Look critically at the layout. Some poster 'experts' suggest that if there is about 20-25% text, 40-45% graphics and 30-40% empty space, you are doing well.



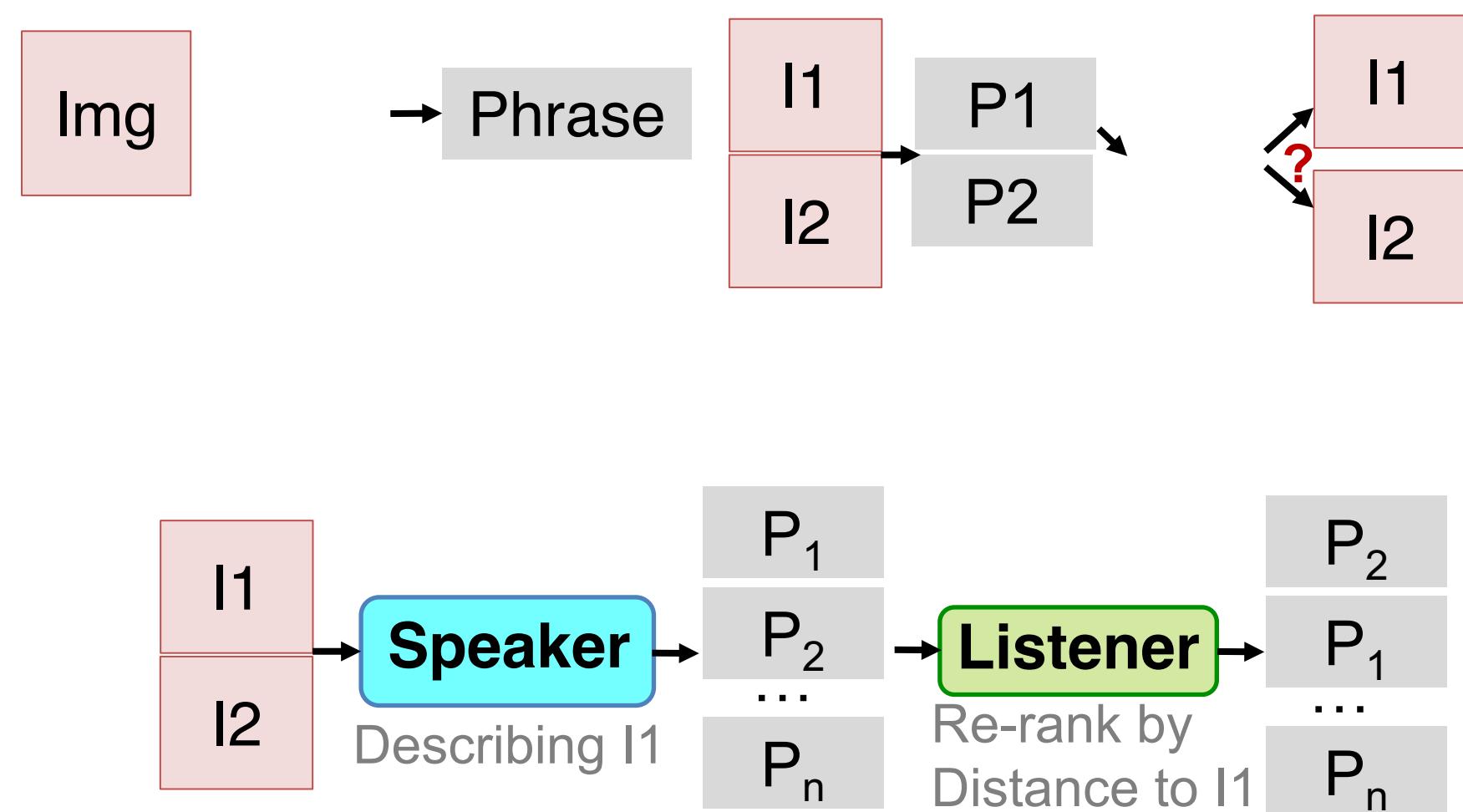
## Section 3:

- Include more figures than are in the paper so you can talk to them. Include things that are not in the paper and then encourage them to read the paper. Don't try to just put all the paper here.
- If it looks like a cut/paste of the paper, people skip that poster since they can read the papers after the conference. Many people find it better to spend time talking with poster presenters that have more to offer than just redoing the paper content paper in big fonts.
- Remember Poster boards look like this.. This is your canvas. Paint us a picture of your work.



Maybe a QRCode to  
the website with your code.

- **SS**: sees and describes one image
- **DS**: sees and describes an image pair



- **SL<sub>r</sub>**: trained on random negative images
- **SL**: trained on original image pairs as collecting the phrases
- **DL**: sees paired phrases

#### Embedding

We show the 1024-dimensional feature of the 500 most frequent phrases, projected into two dimensions using t-SNE.

#### Pragmatic speaker

We use our speaker to generate several phrases, and re-rank them with our listener model.