



Reasoning about Fine-grained Attribute Phrases using Reference Games

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

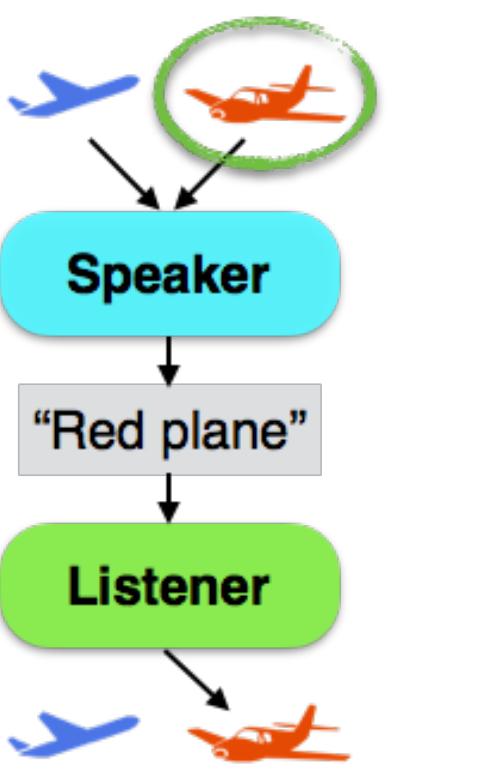
Motivation

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



We collect a dataset based on OID^[1] dataset, containing 9,400 image pairs and 5 descriptions per pair.

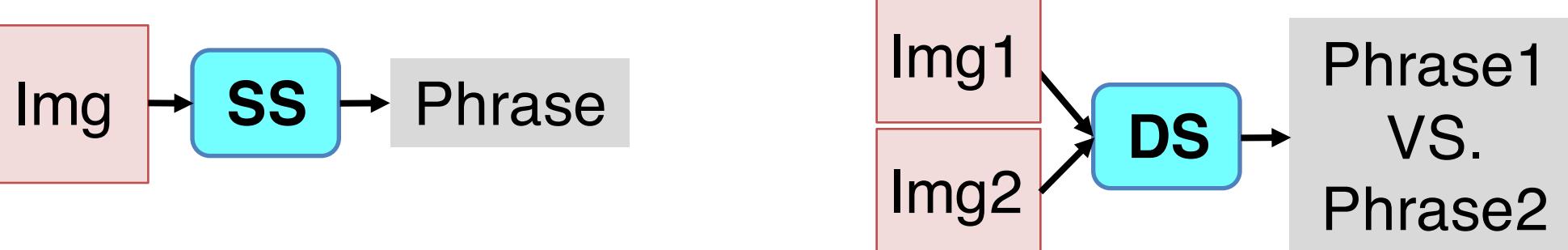
Reference Game: The speaker model describes the attribute of an image. The listener model guesses which images it is referring to.



Approach

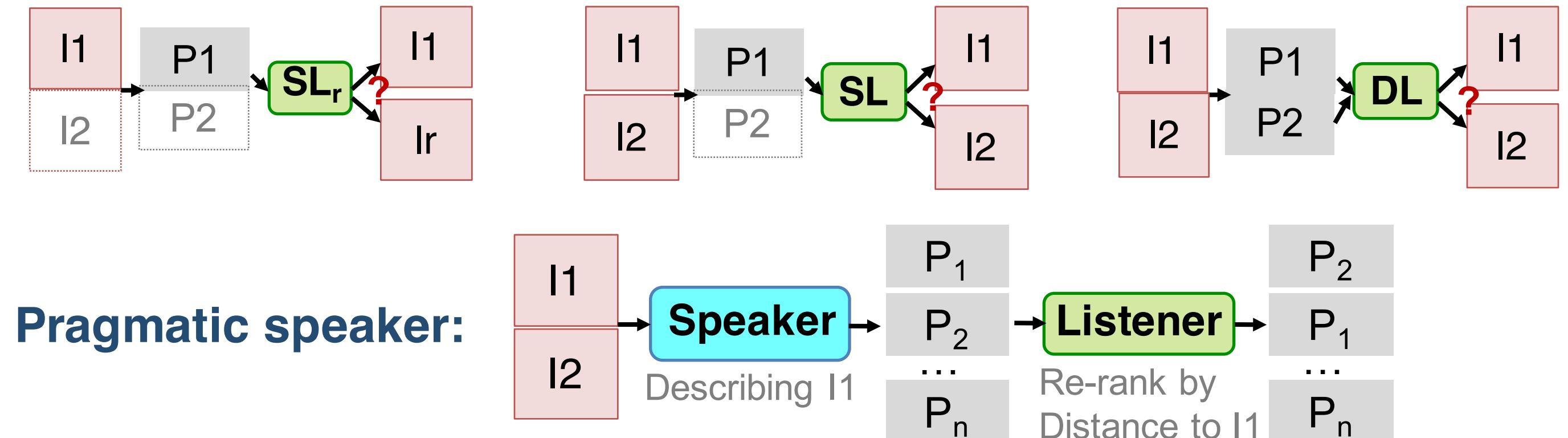
Speaker:

We follow the Show-and-tell model^[2]: a LSTM image captioning model conditioned on vgg-16 image features. [SS, DS]:



Listener:

We train a joint embedding model for images and descriptions with vgg-16 and LSTM, with triplet loss on top. [SL_r, SL, DL]:



Pragmatic speaker:

Quantitative Result

→: Listener

The accuracy of picking the target image out of a pair

| Top | Accuracy (%) | | | Human Test* |
|-----|-----------------------|---------|-------------|-------------|
| | SL _r Test* | SL Test | Human Test* | |
| 1 | 84.0 | 79.8 | 83.0 | 81.7 |
| SS | 80.0 | 79.2 | 78.0 | 80.6 |
| 10 | 78.0 | 78.9 | 76.6 | 80.0 |
| | | | | 64.2 (74.1) |
| 1 | 94.0 | 92.8 | 92.0 | 92.8 |
| DS | 91.2 | 90.3 | 91.2 | 91.4 |
| 10 | 88.6 | 88.8 | 90.0 | 90.5 |
| | | | | 82.0 (88.5) |
| | | | | 80.2 (86.7) |
| | | | | 77.9 (85.0) |

| Input | Speaker | Listener | Val | Test |
|-----------------------------------|---------|-----------------|------|------|
| P ₁ | Human | SL _r | 82.7 | 84.2 |
| | | SL | 85.3 | 86.3 |
| P ₁ vs. P ₂ | Human | DL | 88.7 | 88.9 |
| | | 2×SL | 89.6 | 89.3 |

| Top | Human listener accuracy (%) | | | Reranker listener SL |
|-----|-----------------------------|-----------------|-------------|----------------------|
| | None | SL _r | SL | |
| 1 | 68.0 (77.0) | 94.0 (96.0) | 87.0 (92.0) | |
| SS | 64.2 (74.1) | 82.6 (88.3) | 80.8 (87.1) | |
| 7 | 63.1 (72.8) | 74.3 (82.0) | 74.3 (82.4) | |
| 1 | 82.0 (88.5) | 95.0 (96.5) | 95.0 (97.0) | |
| DS | 80.2 (86.7) | 90.0 (93.3) | 88.6 (92.8) | |
| 7 | 79.1 (85.6) | 86.7 (91.5) | 86.1 (91.1) | |

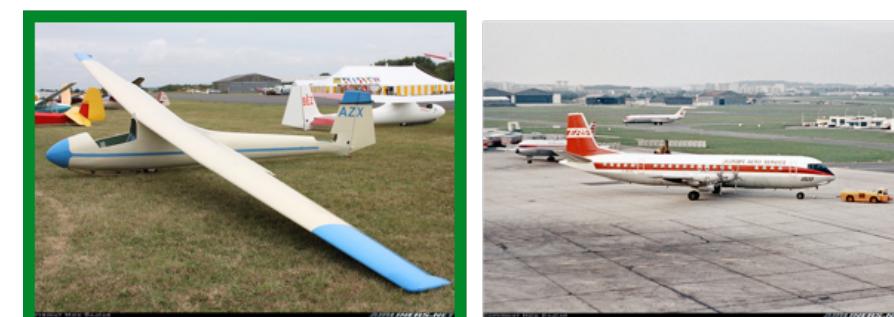
↑: Speaker

Use speaker models to generate phrases, report the accuracy from different listeners

Example Output

→: Speaker

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



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:

- 1) no engine (p=0.2963)
- 2) small (p=0.1800)
- 3) private plane (p=0.0650)
- 4) on the ground (p=0.0519)
- 5) propeller engine (p=0.0322)
- 6) on ground (p=0.0250)
- 7) glider (p=0.0228)
- 8) white color (p=0.0163)
- 9) small plane (p=0.0151)
- 10) no propeller (p=0.0124)

↓: Pragmatic speaker
Re-ranking improves the result!



ss:
✓ passenger plane
? white
✓ jet engine
? facing right
✓ commercial plane
? _UNK
? on the ground
✓ large
✓ large size
? facing right
✓ on runway
✓ on concrete
✓ t tail

ss + SL_r:
✓ commercial plane
✓ facing right
✓ large
✓ large size
✓ jet engine
✓ on concrete
✓ t tail
✓ passenger plane
? on the ground
✓ white and red
✓ white colour with red stripes

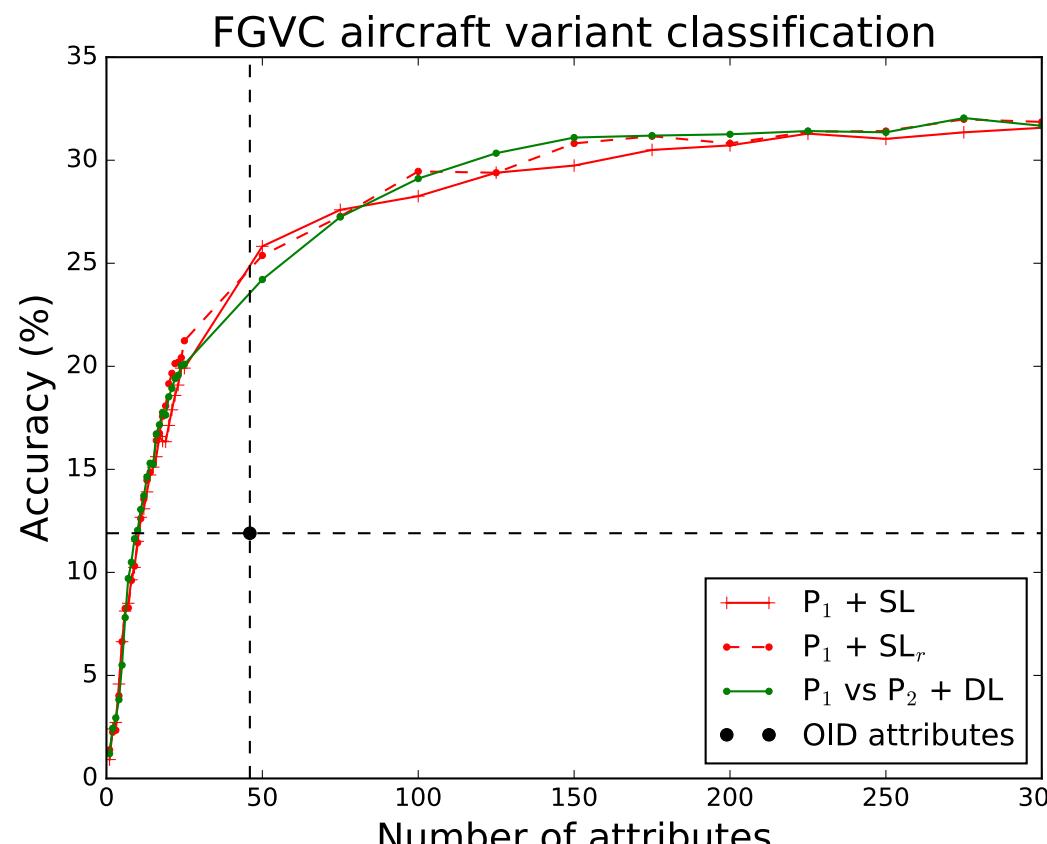
DS:
✓ commercial plane
✓ facing right
✓ turbofan engine
✓ on concrete
✓ t tail
✓ jet engine
✓ twin engine
✓ multi seater
✓ t tail
✓ white and red
✓ white colour with red stripes

DS + SL_r:

- ✓ commercial plane
- ✓ facing right
- ✓ jet engine
- ✓ turbofan engine
- ✓ on concrete
- ✓ t tail
- ✓ jet engine
- ✓ twin engine
- ✓ multi seater
- ✓ t tail
- ✓ white and red
- ✓ white colour with red stripes

Fine-grained Classification

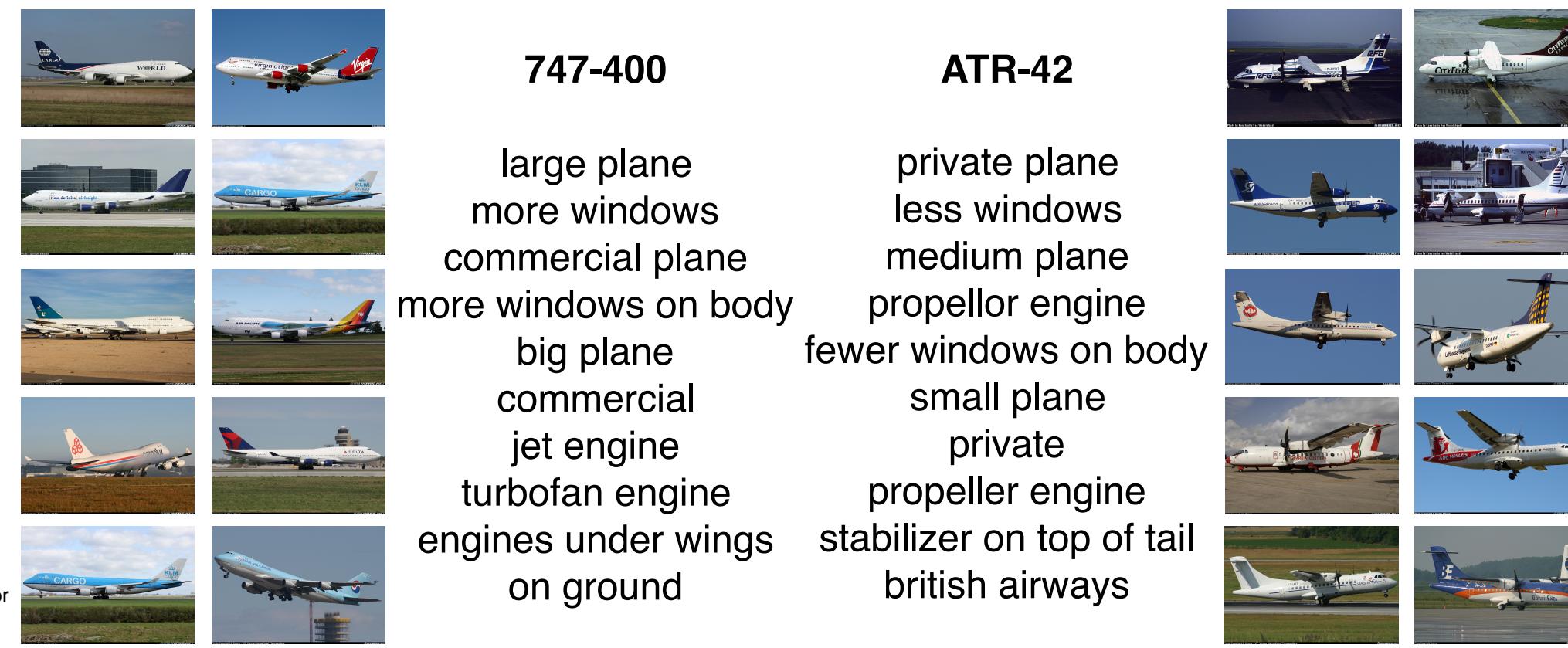
We compare the FGVC-aircraft^[3] classification result using two attributes:
Use our listener model to get scores between the top-k frequent attribute phrases and the image;
Expert-designed 46 attributes from OID^[1] dataset.
Our attribute phrases outperform by **20%**.



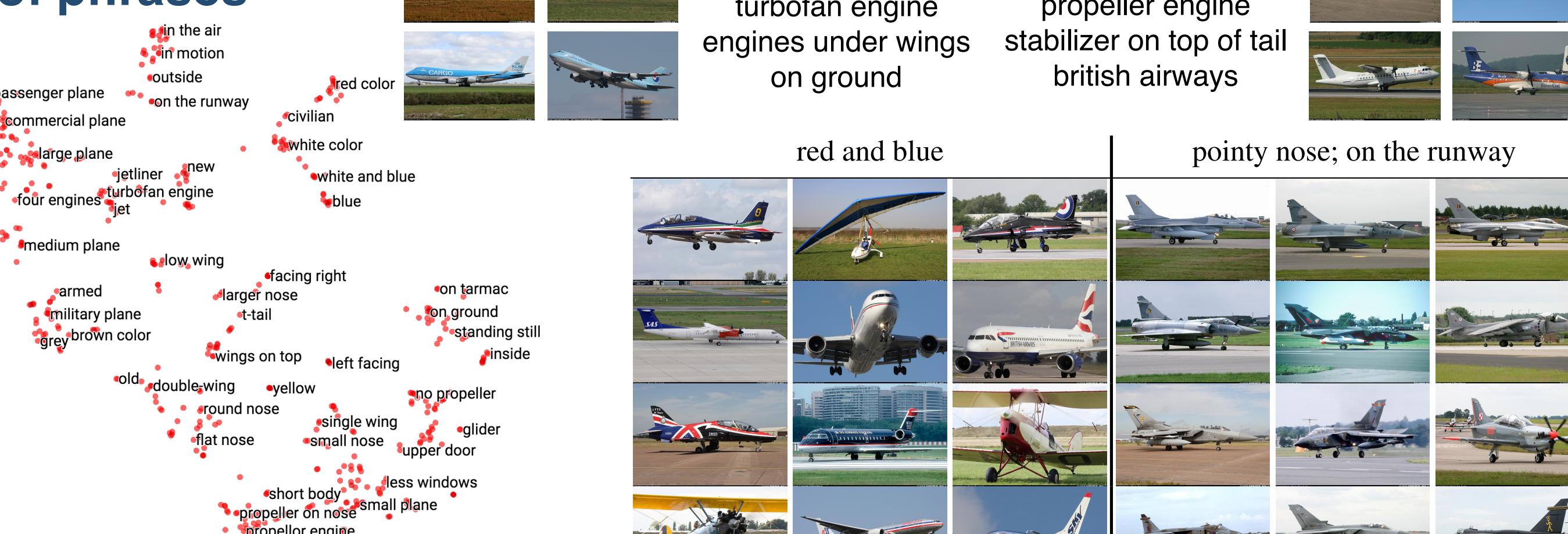
Visualization & Application

→: Set-wise attributes

Vote the speaker output



↓: Embedding of phrases



→: Image retrieval

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



- [1] A. Vedaldi et al., Understanding objects in detail with fine-grained attributes, CVPR, 2014.
[2] O. Vinyals et al., Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge, TPAMI, 2016.
[3] S. Maji et al., Fine-grained visual classification of aircraft. arXiv:1306.5151, 2013.