# Interpretable Machine Teaching via Feature Feedback

Shihan Su

Yuxin Chen Oisin Mac Aodha Pietro Perona

**Yisong Yue** 

Caltech



Thm: The worst-case cost of our greedy strategy achieving error  $\epsilon$  is within a logarithmic factor of the worst-case cost of the optimal algorithm achieving error of at least  $P(h^*)\epsilon/2$ 

# **Results**

## Datasets



Student's ability to learn a new concept can be greatly improved by providing them with clear and interpretable **explanations** from a knowledgeable teacher

• Each with 128 images of two balanced classes: Jupiter and Mars

















#### **Decision Rule**

- Predictive features: parts indicated by arrows
- Hard Mars: blue top, square middle and thick base Jupiter: any other combination
- Easy Mars: yellow top, circle middle and large ellipse near bottom lupiter: any other combination

#### **Baseline**

Random: Random images with no explanations Random-feature: Random images with random explanations STRICT: Label-based greedy approach

### Conclusion

With explanation based machine teaching, students achieve

- Better accuracy
- Faster question answering at test time

