# Purah P2 Lightning Talk

Ankit Lalwani, Shihan Ma, Nathan Wynn

# Ethics - Key Takeaways

| Collection | Data Storage | Analysis | Modeling | Deployment |
|---|---|---|---|---|
| Large upfront cost, large payoff in the modeling stage.<br><br>YFCC100m pros and cons: diversity and skew.<br><br>Creative Commons license! | Data retention plans were unnecessary. | Higher diversity -> lower bias towards physical characteristics.<br><br>Higher bias towards online Flickr community. | Transparency.<br><br>Only using facial anatomic features, no personal information, no categorical information. | Full control of rollbacks, no feature drift.<br><br>Educational purposes only! |

# Python Packages

- Pandas
  - Read CSV files
  - Select useful features (39 features)


- Scikit-learn
  - K-th nearest neighbors (KNN)
  - Linear Discriminative Analysis (LDA)
  - Quadratic Discriminant Analysis (QDA)

# ML models

- Classification
    - K-th Nearest Neighbors (second best accuracy)
    - Linear Discriminative Analysis (best accuracy)
    - Quadratic Discriminant Analysis
    - Random Forest
    - Logistic Regression

- dimensionality reduction
    - Principal Component Analysis
    - Elastic Net Regression (with cross-validation)
        - Selects input variables with high coefficients.

# Thank you!

Questions?