

Introduction To Explainable AI (XAI)

Rohan Mitra

Machine Learning Engineer @ Bayut & Dubizzle

Part-time Research Assistant @ CUD

Organizer & Speaker @ Google Developer Groups Sharjah

Case Study: Mount Sinai Hospital



- Wanted to use AI to detect pneumonia from chest X-Rays
- Built a CNN model
 - Trained on lots of X-Rays from their own hospital
 - 158,323 X-Rays used!
 - Made sure to have a proper train-test split
- 99.95% accuracy!!



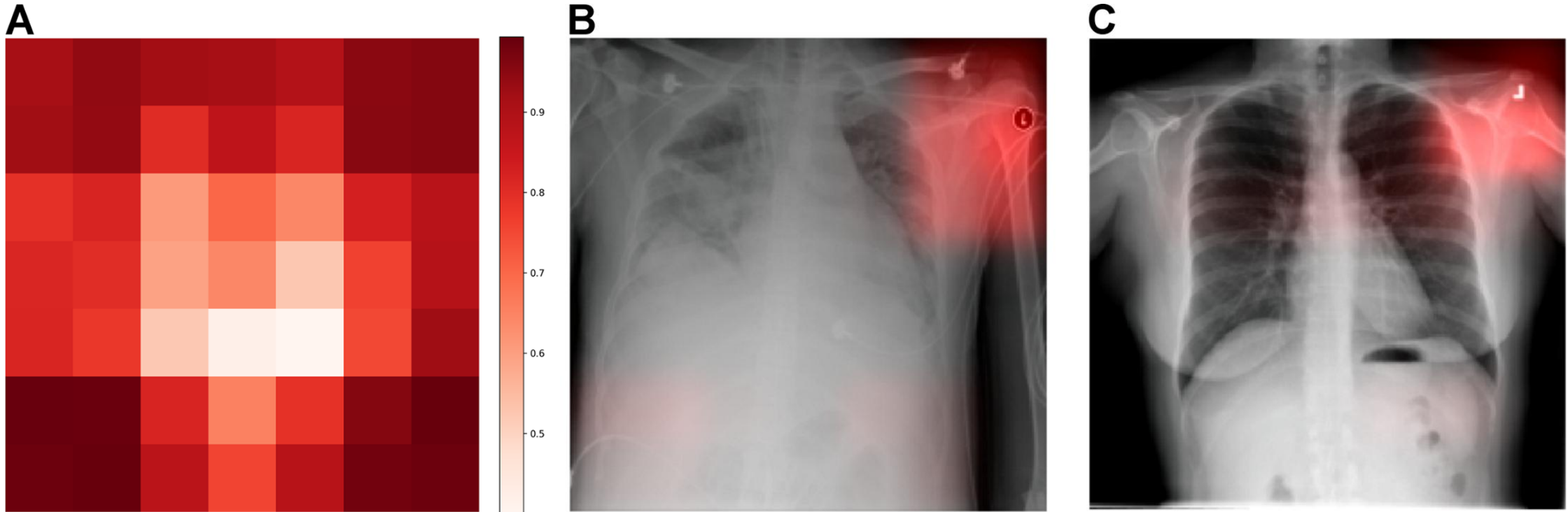
Case Study: Mount Sinai Hospital

- “We want to share this amazing model!”
- Contact National Institutes of Health Clinical Center (NIH)
- They loved it and got it implemented!
- 72% accuracy...

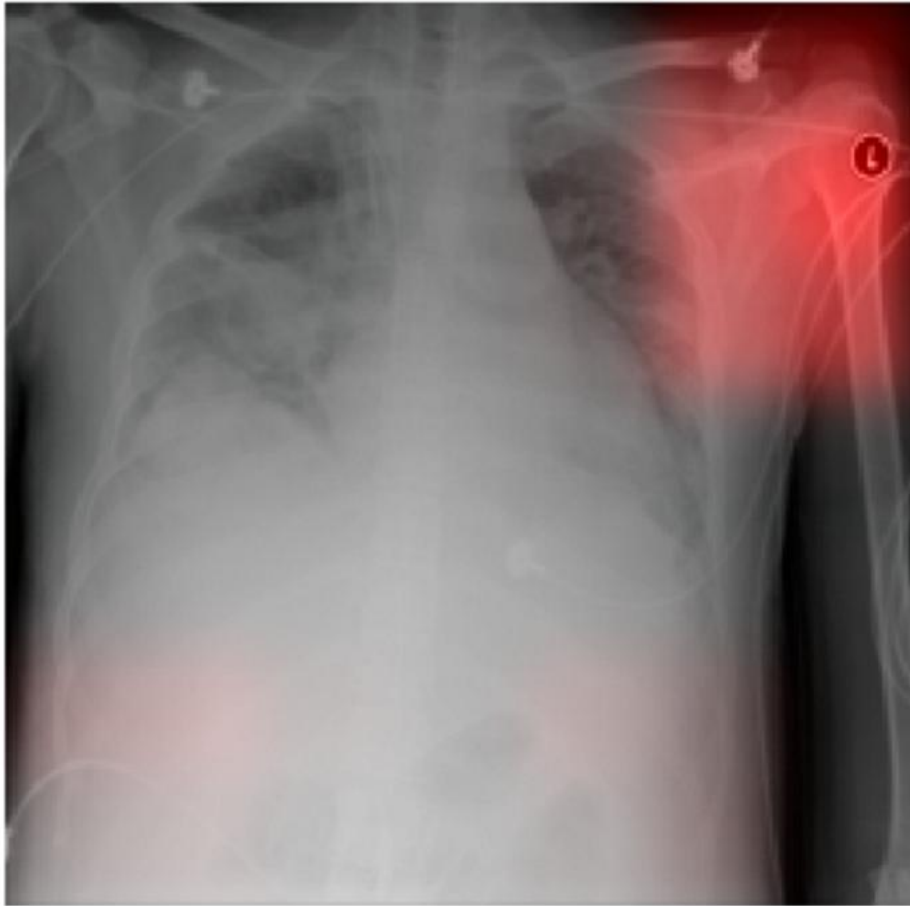


Case Study: Mount Sinai Hospital

- XAI to the rescue!
- Examine activation maps and overlay on the image



Case Study: Mount Sinai Hospital



- Mount Sinai physicians put a metal physical marker on the patients' shoulders!
- The CNN took a shortcut!

Impact of XAI

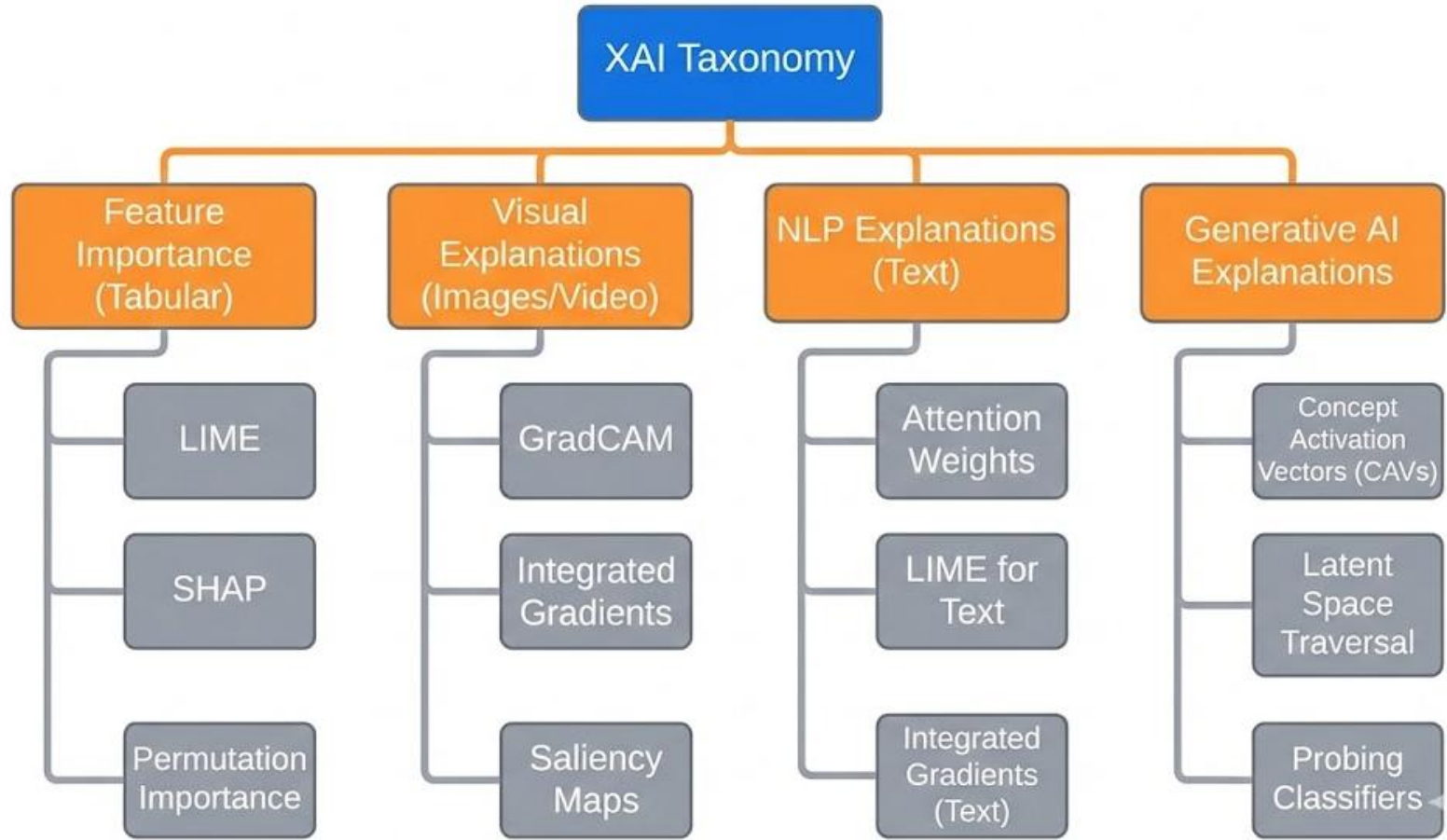
Trust & Transparency (External)

- **Focused on Users:** Doctors, Regulators, Buyers
- Users demand explanations
- High stakes predictions: Medicine, Finance, Law
- Ethical introspection

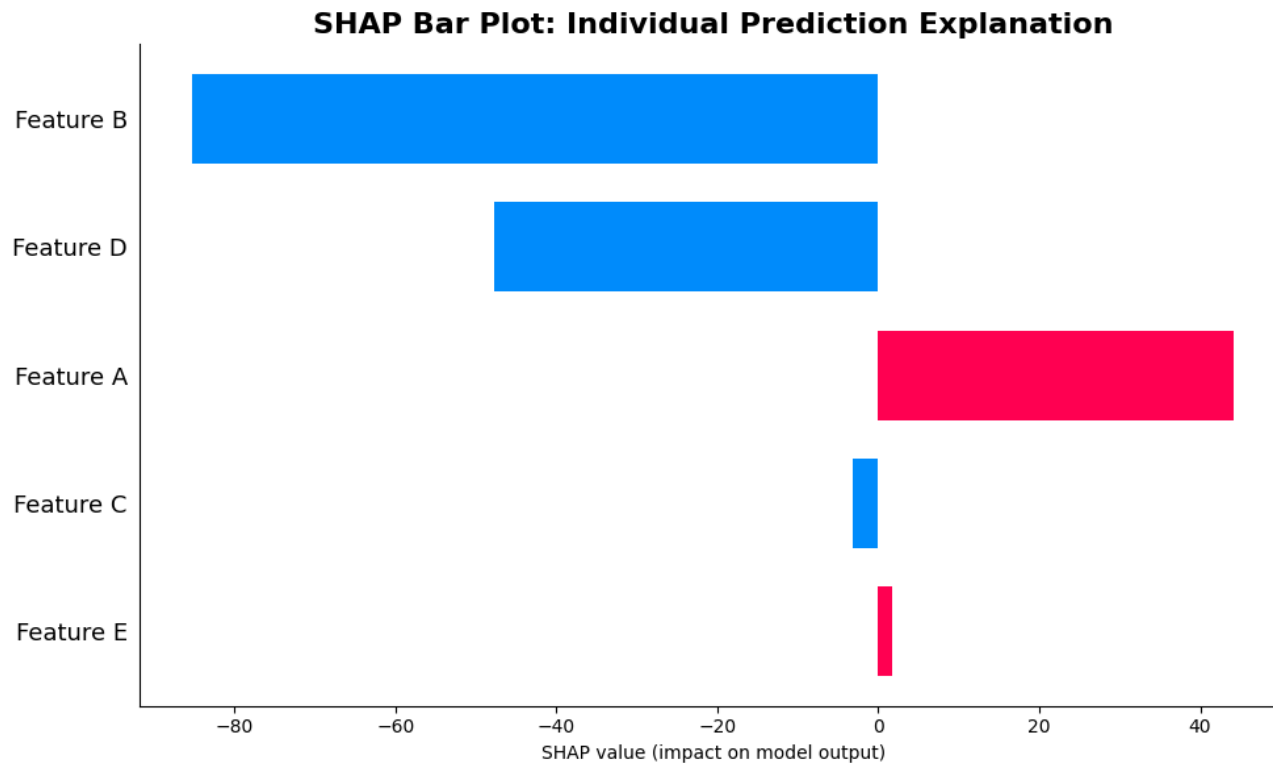
Debugging & Model Hygiene (Internal)

- **Focused on Builders:** Data Scientists, Machine Learning Engineers, YOU!
- Treated as a unit-test to ensure model learns useful features
- Models tend to find shortcuts
- Avoid spurious correlations

Taxonomy



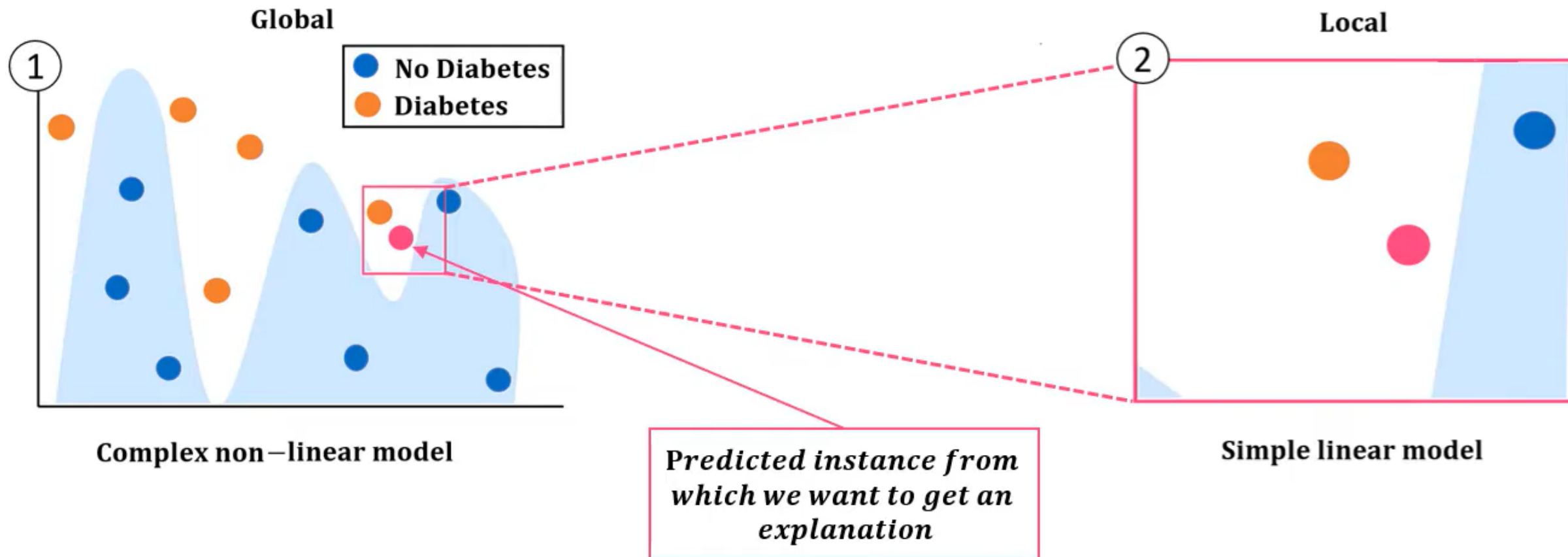
XAI Method: SHAP



- **Origin:** Based on Cooperative Game Theory (Lloyd Shapley, Nobel Prize)
- For a given prediction, much did each feature contribute to the final prediction?
- Calculates the **Marginal Contribution** of a feature across all possible combinations of features
- **Global AND Local:**
 - *Local:* Why did **this** specific person get denied a loan?
 - *Global:* What drives the model generally?

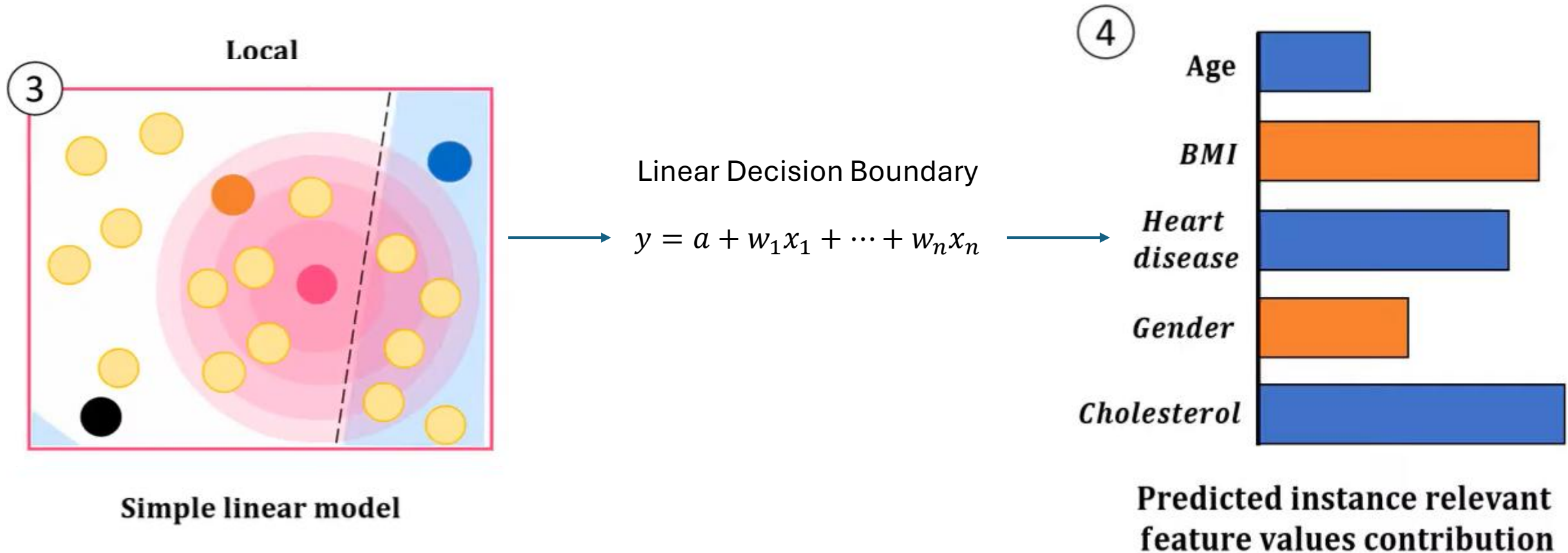
XAI Method: LIME

- Local Interpretable **M**odel-agnostic **E**xplanations
- Approximate the global decision boundary with linear approximations near a point
- We don't need to understand how the entire model works, just how it behaves near this particular prediction



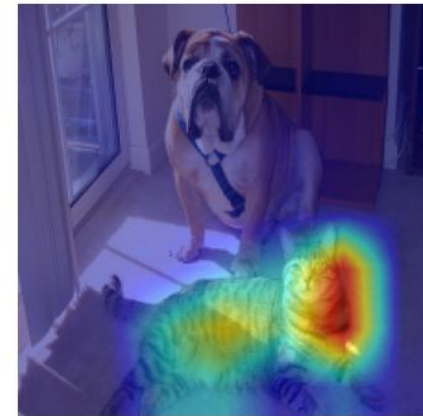
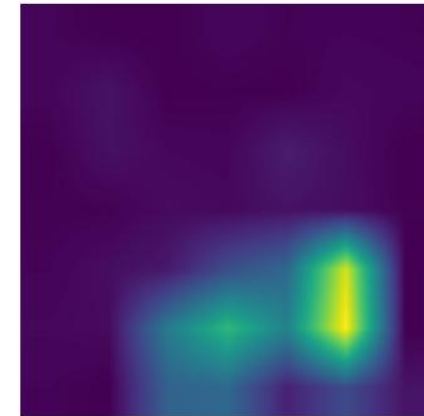
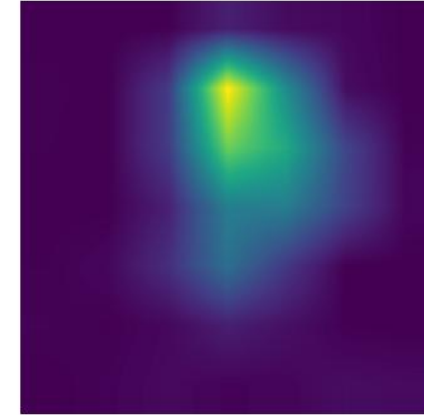
XAI Method: LIME

- Local Interpretable Model-agnostic Explanations
- Approximate the global decision boundary with linear approximations near a point
- We don't need to understand how the entire model works, just how it behaves near this particular prediction

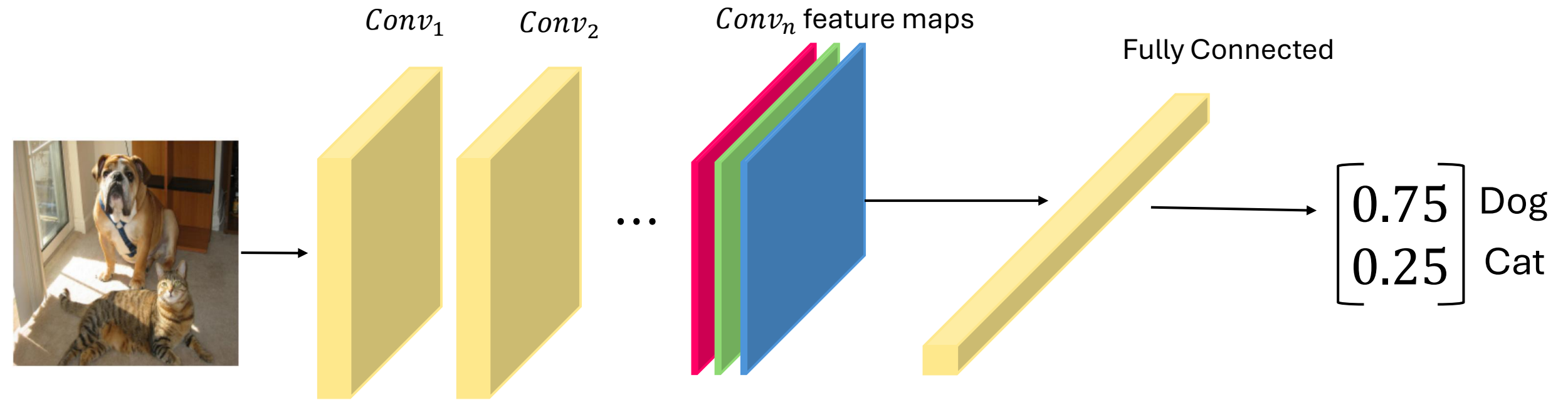


XAI Method: Grad-CAM

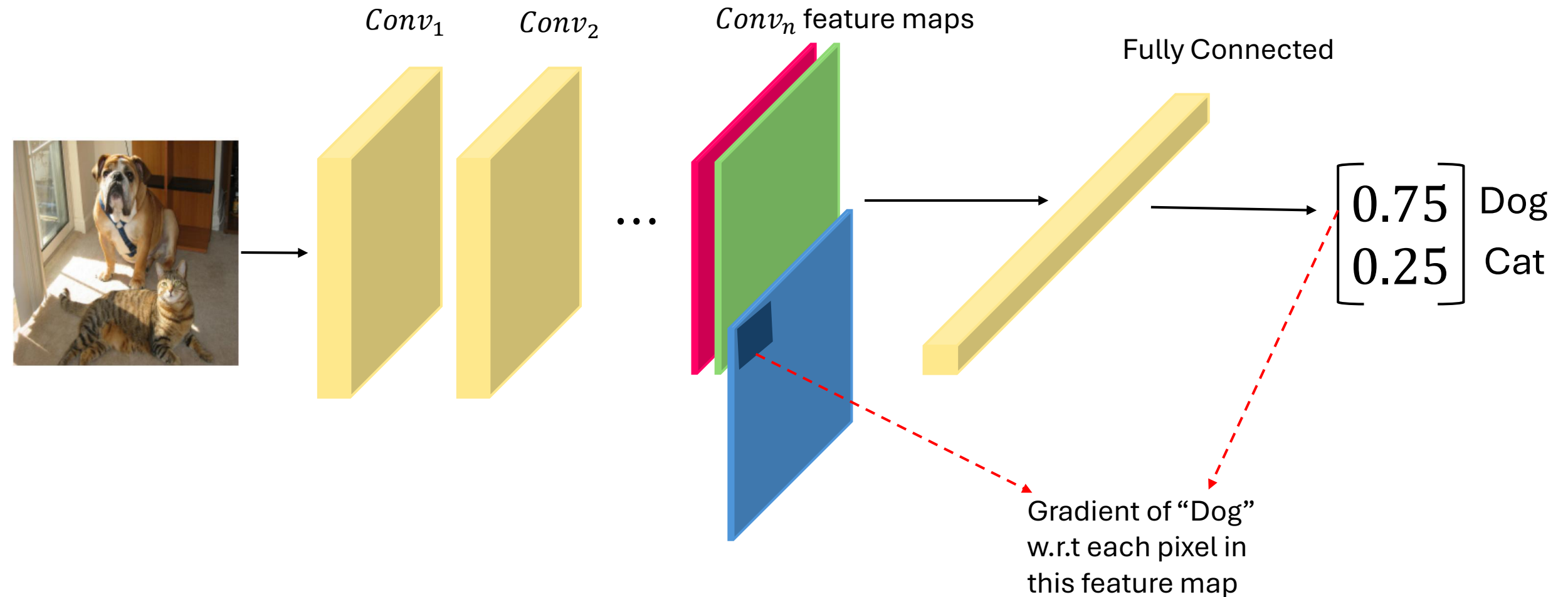
- **Gradient-weighted Class Activation Mapping**
- For image/video data
- Focuses on inner workings of the models
 - Looks at how the final class probability changes with respect to the features changing
- Generates a heatmap of the areas in the image that were important to determine the class



XAI Method: Grad-CAM

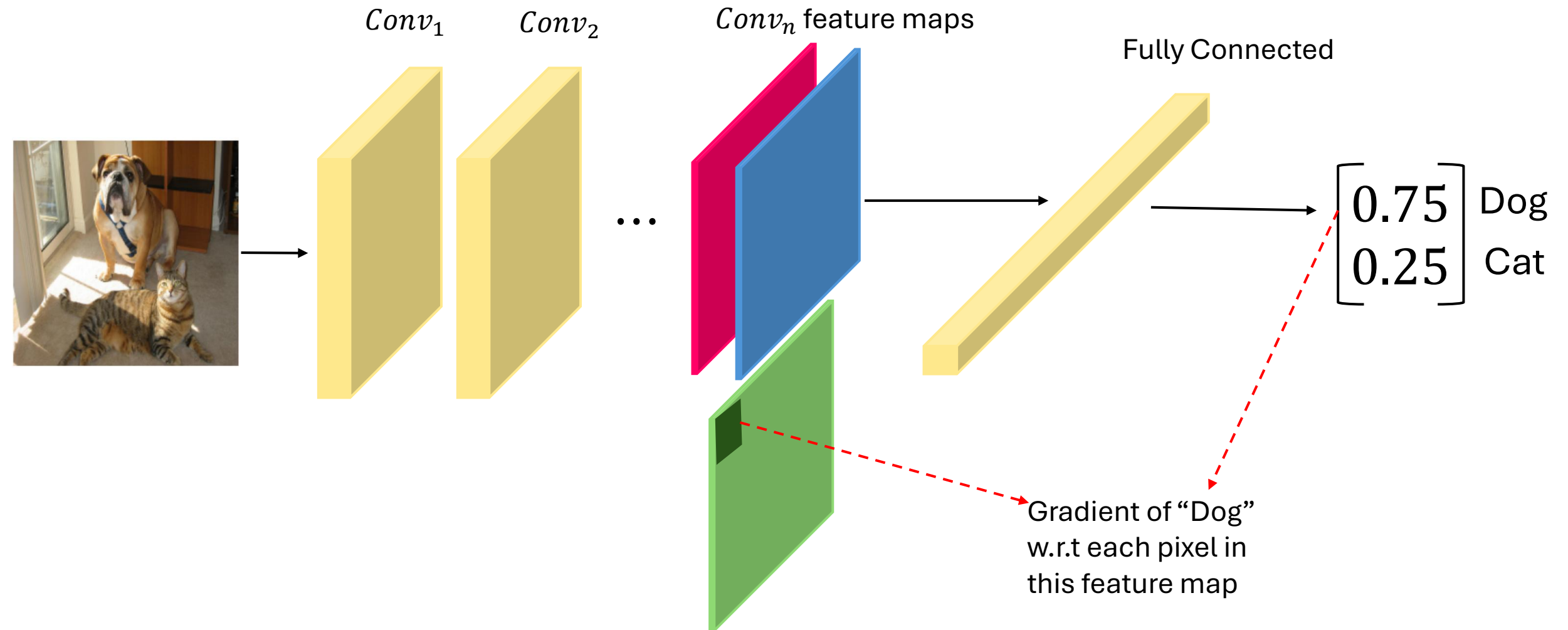


XAI Method: Grad-CAM

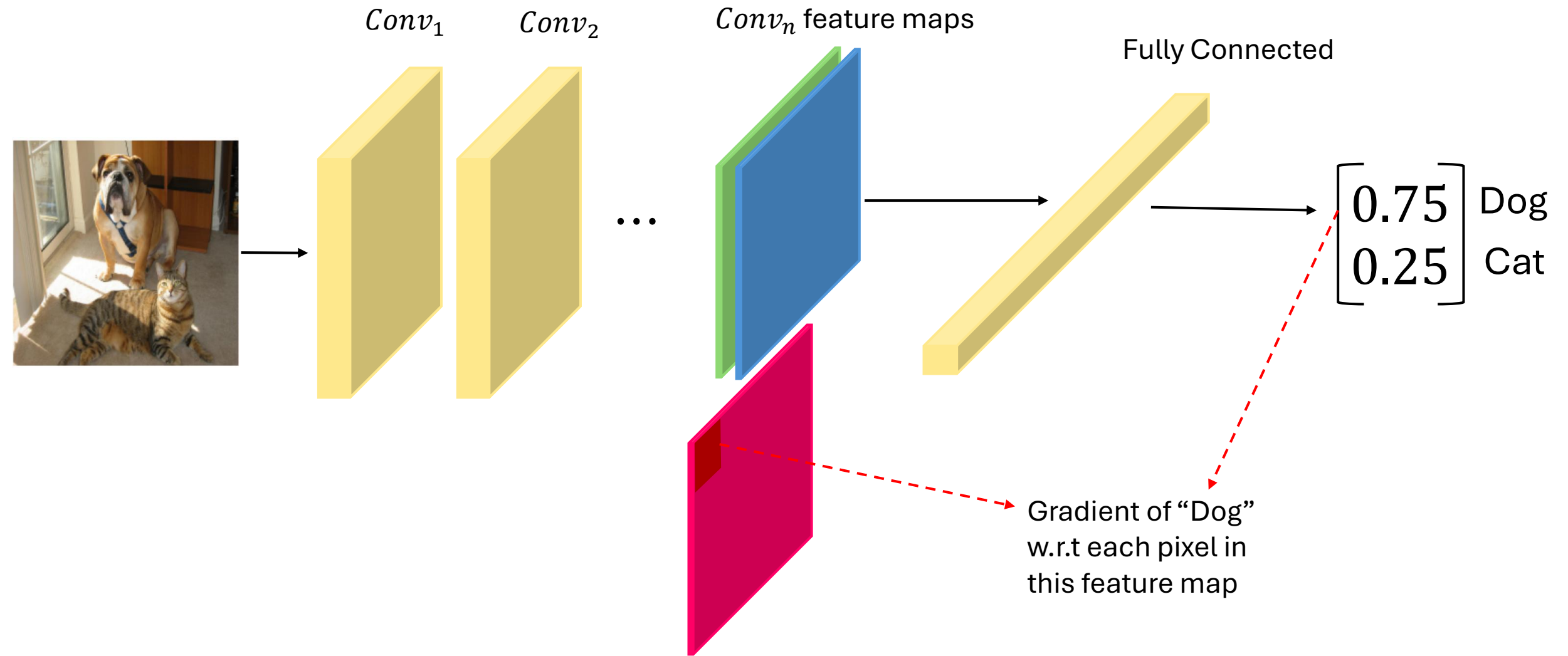


How much does the probability of "Dog" change when this pixel changes?

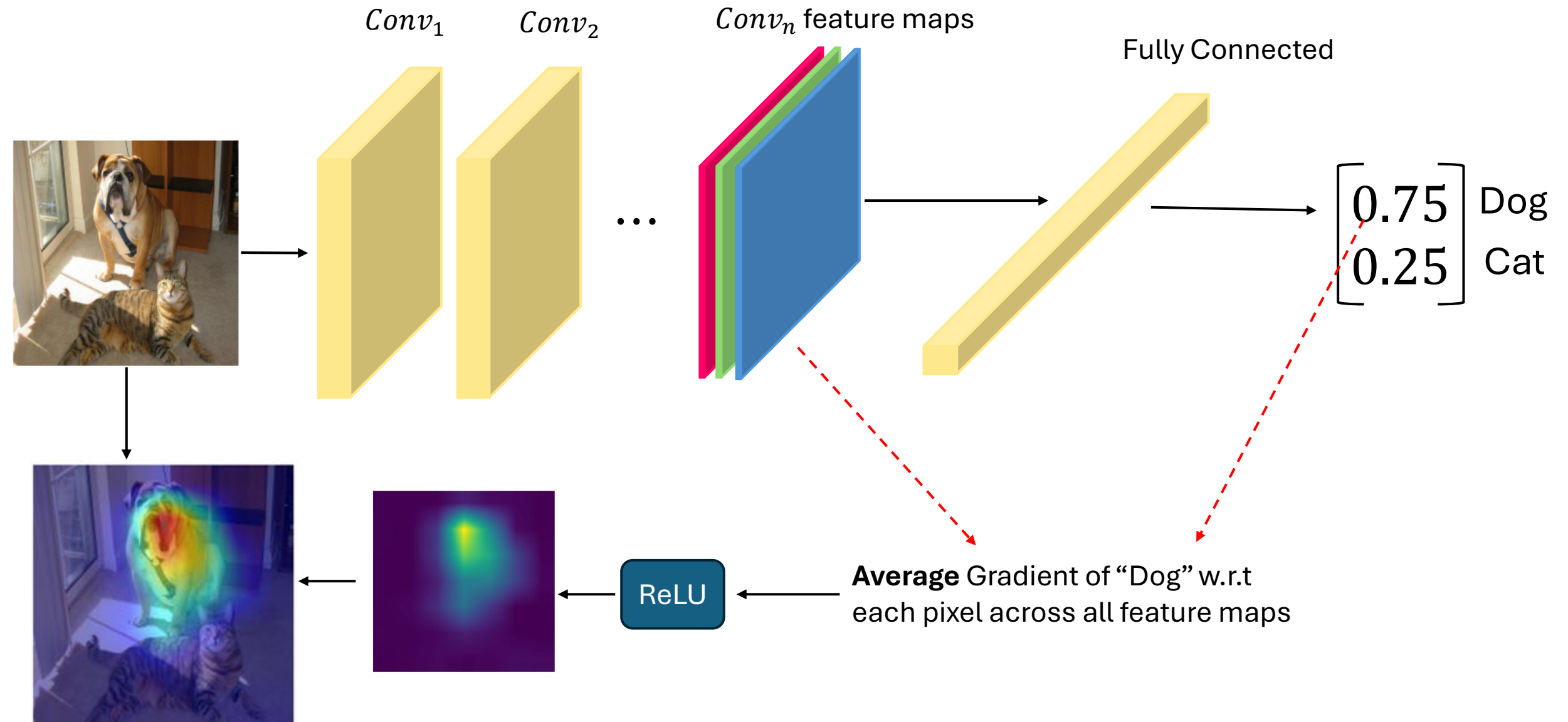
XAI Method: Grad-CAM



XAI Method: Grad-CAM



XAI Method: Grad-CAM





Let's see some code!

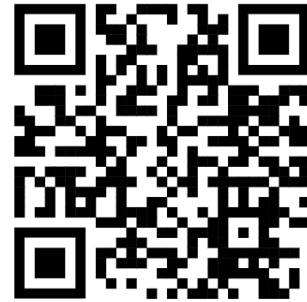


github.com/ro1406/xai-tutorials

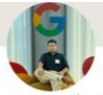


Thank You!

My Website



rohanmitra.dev



Rohan Mitra

Machine Learning Engineer @Bayut & Dubizzle
(OLX Group) | Research Assistant for Machine Lea...



linkedin.com/in/rohan-mitra14/