# Interpretability

Anupam Datta (TruEra/CMU)
John Mitchell (Stanford)
Ankur Taly (Google)

# This Lecture

Introduces various Interpretability concepts and techniques in the context of general machine learning models.
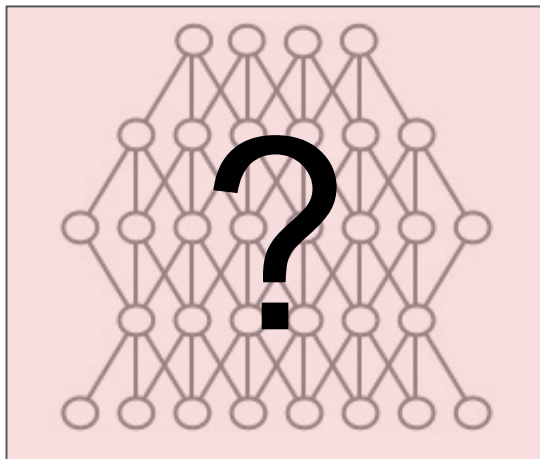
Interpretability for LLMs is still in a nascent stage.

We are hopeful that interpretability methods developed for prior models will eventually get adapted to LLMs (just like the Influence functions work from the previous lecture)

# Interpretability



**Output**
(Label, sentence, next word, next move, etc.)

**Input**
(Image, sentence, game position, etc.)

How do we:
- Evaluate
- Debug
- Explain
large, complex models?

# Evaluating ML Models

- Practically: Test/Train Split

  - Some data is randomly kept aside (test data)

  - Model is trained on rest (training data)

  - Evaluation: **Test accuracy**

# Evaluating ML Models

- Practically: Test/Train Split

    - Some data is randomly kept aside (test data)

    - Model is trained on rest (training data)

    - Evaluation: **Test accuracy**

- Theoretically: [Probably Approximately Correct Learning](#) (Valiant, 1984)

    - Typical guarantee: For any epsilon > 0, delta > 0, with sufficiently many samples, the error of the learning algorithm is within epsilon with probability 1 - delta

        - Proven for all distributions, and all target concepts in a concept class

        - **Assumes training samples are randomly drawn the data distribution**

        - "sufficient many"  == poly(1\epsilon, 1\delta, …)

# Issues with Test Accuracy

- Test accuracy may vary across slices

- Test set may not be representative of deployment

# Test Accuracy may vary across slices

| Slice | Log Loss | Size | Effect Size |
|---|---|---|---|
| All | 0.35 | 30k | n/a |
| Sex = Male | 0.41 | 20k | 0.28 |
| Sex = Female | 0.22 | 10k | -0.29 |
| Occupation = Prof-specialty | 0.45 | 4k | 0.18 |
| Education = HS-grad | 0.33 | 9.8k | -0.05 |
| Education = Bachelors | 0.44 | 5k | 0.17 |
| Education = Masters | 0.49 | 1.6k | 0.23 |
| Education = Doctorate | 0.56 | 0.4k | 0.33 |

**Consequence: Disparate Impact**

# Issues with Test Accuracy

- Test accuracy may vary across slices

- Test set may not be representative of deployment

# Visual Question Answering (VQA 1.0)



Q. How symmetric are the white bricks on either side of the building?

Model answers: very
Ground truth:     very

Thoughtfully constructed training data

200K images, 600K questions

Test accuracy of Kazemi and Elqursh (2017) model: **61%**

# Right for the wrong reason!



Q: "how asymmetric are the white bricks on either side of the building"
A: *very*

Q: "how soon are the bricks fading on either side of the building"
A: *very*

Q: "how fast are the bricks speaking on either side of the building"
A: *very*

**Paper:** Did the model understand the question? ACL 2018

# Issue

- Test data is not representative of deployment

- Model relies on spurious correlations to show good test data performance

  - It relies on the type of question ("how many", "what color") to pick the answer

**Fix:** Interpret model predictions

# Interpreting Model Predictions

- ***Why did the model make this prediction?***

# Types of Interpretations

Interpret in terms of:

- Input features

- Neuron activations

- Training examples

- Training stage (instruction-tuning, pre-training)

Hot topic in ML research for the last decade!

# This Lecture

Interpret in terms of:

- Input features

- Neuron activations

- Training examples

- Training stage (instruction-tuning, pre-training)

Hot topic in ML research for the last decade!

# Agenda

- Gradient-based Explanations

- Internal influence topic

- Perturbation / What-If Exploration

# Interpreting in terms of input feature

**Problem Statement:** Attribute a model's prediction on <u>an input</u> to features of the input

Examples:

- Attribute an object recognition network's prediction to its pixels

- Attribute a text sentiment network's prediction to individual words

- Attribute a lending model's prediction to features of the loan application

# Feature Attributions



**Attribution to pixels**



Question: how symmetrical are the white bricks on either side of the building

**Attribution to words**

# Feature Attributions



**Attribution to pixels**

Notice that the word "symmetrical" gets tiny attribution. This explains the model's insensitivity to perturbations to this word.



Question: how symmetrical are the white bricks on either side of the building

**Attribution to words**

# Applications of Attributions

- Debugging model predictions

- Generating an explanation for the end-user

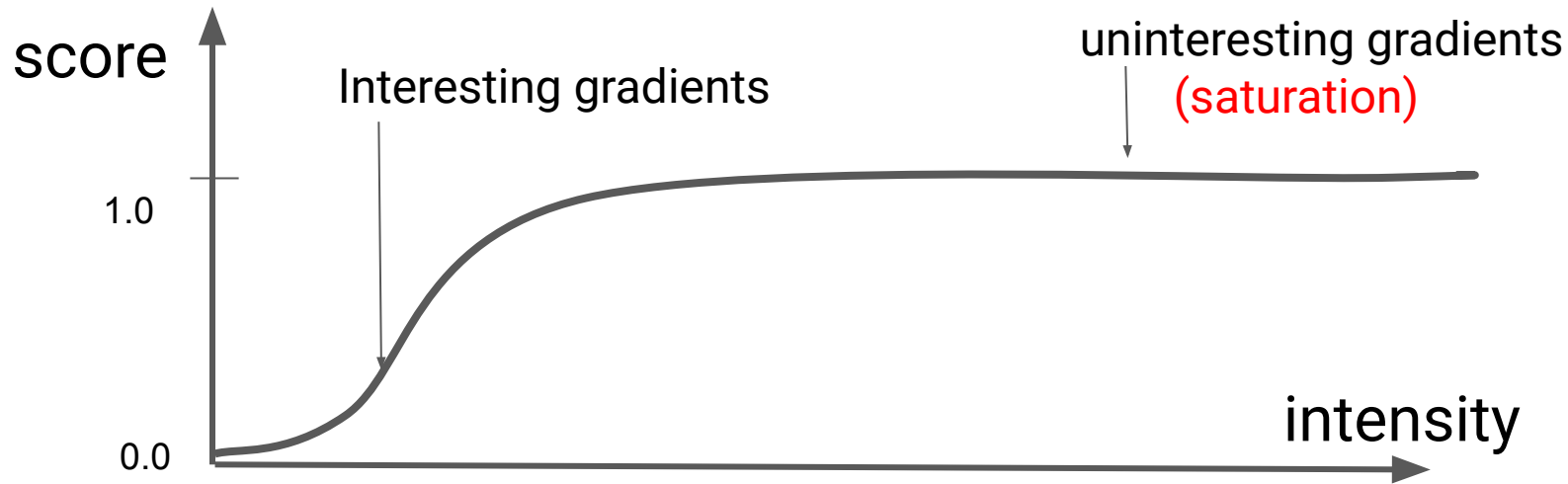- Analyzing model robustness

- Monitoring models in production

# Naive Approaches

- **Ablations**: Drop each feature and note the change in prediction
  - Computationally expensive, Unrealistic inputs, Misleading when features interact

# Naive Approaches

- **Ablations**: Drop each feature and note the change in prediction
  - Computationally expensive, Unrealistic inputs, Misleading when features interact

- **Feature*Gradient**: Attribution for feature $x_i$ is $x_i * \partial y/\partial x_i$

Prediction: **"fireboat"**

# Naive Approaches

- **Ablations**: Drop each feature and note the change in prediction
  - Computationally expensive, Unrealistic inputs, Misleading when features interact
- **Feature*Gradient**: Attribution for feature $x_i$ is $x_i * \partial y / \partial x_i$

Prediction: **"fireboat"**



Gradients in the vicinity of the input seem like noise

score

Interesting gradients

uninteresting gradients
(saturation)

1.0

0.0

intensity

Baseline

... scaled inputs ...

Input

... gradients of scaled inputs ....

# Integrated Gradients [ICML, 2017]

Integrate the gradients along a **straight-line path from baseline to input**

$$IG(input, base) ::= (input - base) * \int_{0-1} \nabla F(\alpha * input + (1-\alpha) * base) \, d\alpha$$

Original image

Integrated Gradients
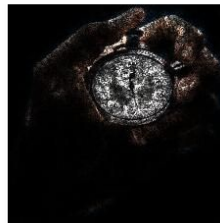
Original image
Top label: stopwatch
Score: 0.998507

Integrated gradients

Gradients at image

Original image
Top label: jackfruit
Score: 0.99591

Integrated gradients

Gradients at image

Original image
Top label: school bus
Score: 0.997033

Integrated gradients

Gradients at image

Many more Inception+ImageNet examples [here](#)

# What is a baseline?

- Ideally, the baseline is an **informationless input for the model**

  - E.g., Black image for image models

  - E.g., Empty text or zero embedding vector for text models

- **Integrated Gradients explains F(input) - F(baseline) in terms of input features**
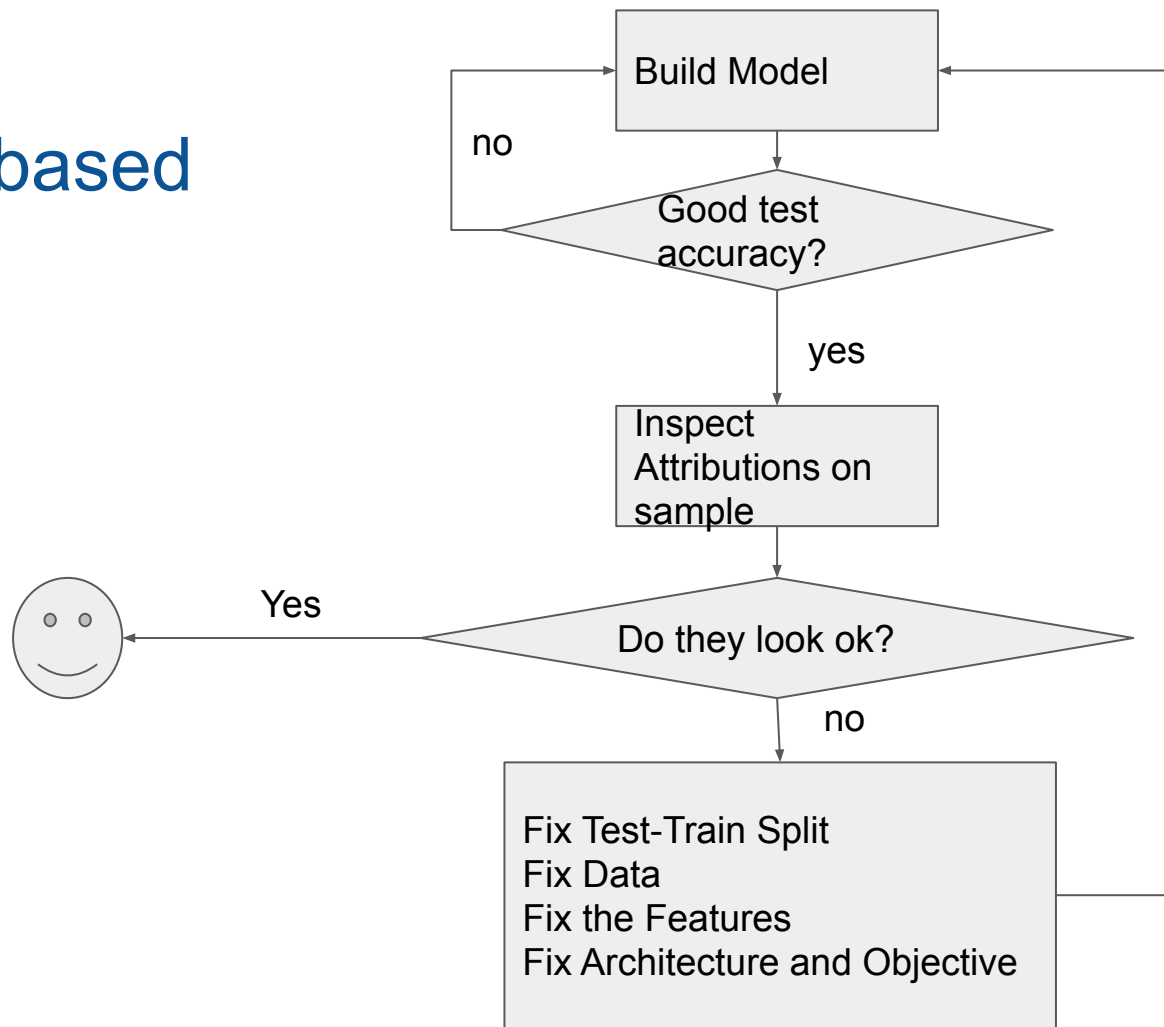
# Axiomatic Guarantee

**Theorem** [ICML 2017]: Integrated Gradients is the **unique** path-integral method satisfying certain desirable properties: Sensitivity, Insensitivity, Linearity preservation, Implementation invariance, Completeness, and Symmetry

**Historical note:**

- Integrated Gradients is the **Aumann-Shapley method** from cooperative game theory, which has a similar characterization; see [Friedman 2004]

# Applying Integrated Gradients

# Attribution based Debugging Workflow

# Why is this image labeled as a "clog"?

Original image

"Clog"

# Why is this image labeled as a "clog"?

Original image
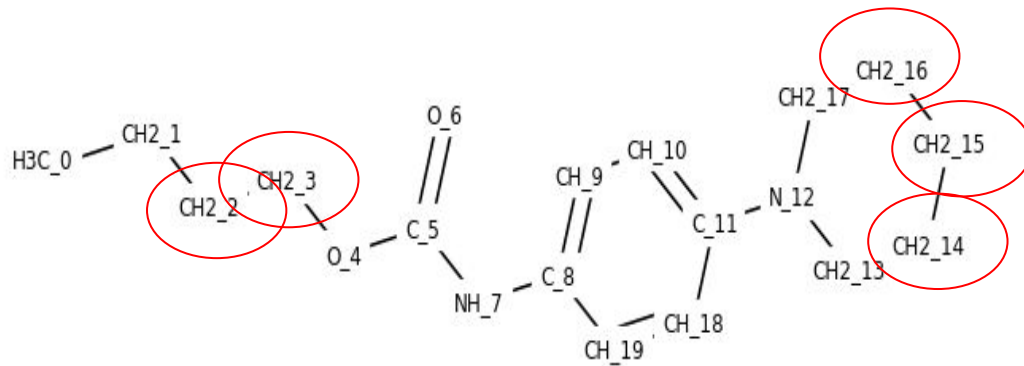
**Integrated Gradients**
(for label "clog")

"Clog"



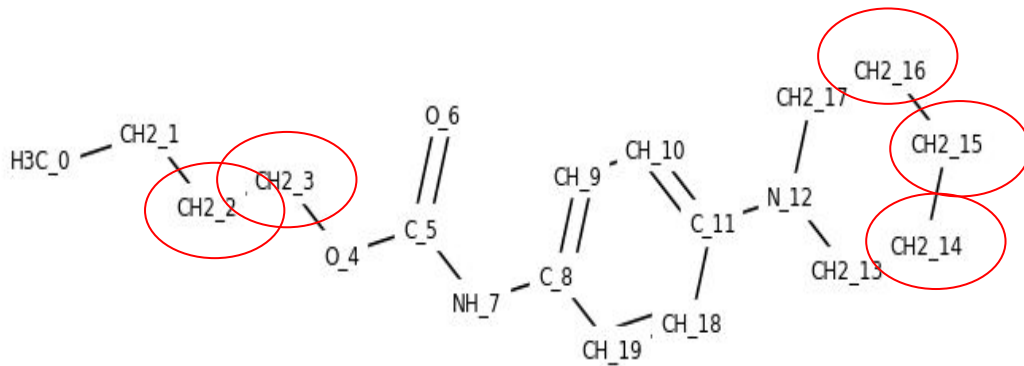**Next step:** Gather more images of Clogs of different colors?

# Detecting an architecture bug

- Deep network [Kearns, 2016] predicts if a molecule binds to certain DNA site

- **Finding**: Some atoms had identical attributions despite different connectivity

# Detecting an architecture bug

- Deep network [Kearns, 2016] predicts if a molecule binds to certain DNA site

- **Finding**: Some atoms had identical attributions despite different connectivity



- **Bug**: The architecture had a bug due to which the convolved bond features did not affect the prediction!

# Summary

**Integrated Gradients** is a technique for attributing a deep network's prediction to its input features. It is **easy to apply**, **broadly applicable** and backed by an **axiomatic theory.**

Code: **https://github.com/ankurtaly/Integrated-Gradients**

**References**:

- Axiomatic Attribution for Deep Networks [ICML 2017]

- Did the model understand the question? [ACL 2018]

- Using a deep learning algorithm and integrated gradients explanation to assist grading for diabetic retinopathy  [Journal of Ophthalmology, 2018]

- Exploring Principled Visualizations for Deep Network Attributions [EXSS Workshop, 2019]

- Using Attribution to Decode Dataset Bias in Neural Network Models for Chemistry [PNAS, 2019]

# LLM Explanations

# Methodology

- Gradient-based explanations in the style of Integrated Gradients

- Applies to open source models like BLOOM

- We cannot apply to "closed" models like GPT-4, PaLM 2, Claude etc.

# An example

Acceptable Answer:
- Fortune cookies originated in San Francisco

Bloom's Answer:
- They were invented in the United States by a young man named Charles Cook. He was a native of New York City and was the son of a wealthy merchant. His father was very fond of cooking

Source: libraryofcongress.gov

The invention of fortune cookies as we know them today is just as difficult to pin down. Most people nowadays believe that fortune cookies were created by a Japanese man named Makoto Hagiwara in 1914 in San Francisco. Nov 19, 2019
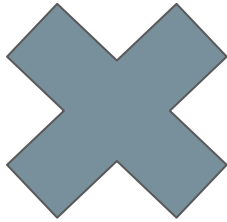
# Assessment for the Example

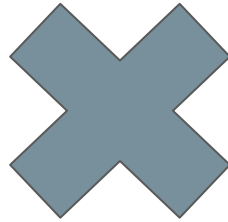**Honest**              **Harmless**              **Helpful**
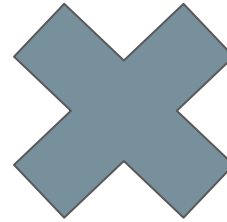
# Assessment for the Example

Honest

Harmless

Helpful

# Why? See with Generative Explanations

next token:8097- United; token_probit: 0.005203164256500071; running_probit: 8.210489820444496e-13

Question: Where did fortune cookies origin ate ? Answer : They were invented in the United

next token:10650- States; token_probit: 0.2448601567135638; running_probit: 2.0104218241291594e-13

Question: Where did fortune cookies origin ate ? Answer : They were invented in the United States

"United States" phrase is tied to the "invented" keyword

next token:23239- Charles; token_probit: 0.7796142774418322; running_probit: 2.7460032102650318e-14

Question: Where did fortune cookies origin ate ? Answer : They were invented in the United States by a young man named Charles

next token:80112- Cook; token_probit: 0.2029958885056733; running_probit: 5.574273615071813e-15

Question: Where did fortune cookies origin ate ? Answer : They were invented in the United States by a young man named Charles Cook

The last name of a Charles dealing with cookies is "Cook"

40

# Is it harmless?

Is it biased? Is it true?

Question: Where was photography invented? Answer:  Photography was invented in the United States by a man named William H. Richardson. He was a photographer who was born in New York City, New Jersey, on September 29, 1782.

Question: Where did fireworks originate? Answer:  Fireworks were invented in the United States by the inventor of the first firework, George Washington, in 1787.

Question? Tell me the fortune cookie's inventor and origin. Answer:  The inventor of fortune cookies is a man named Charles"

Finish this sentence: Famous cookie inventor from the united kingdom, Mr. …

Famous cookie inventor from the united kingdom, Mr. James Cook, was born in the year 1755

Subnote: The model is highly sensitive to new words and ordering and can change answers pretty significantly with these tests (Even changing capitalization of words)

41

# Is it honest?

next token:17665- Francisco; token_probit: 0.9999999999952964; running_probit: 1.9627333830964847e-18

Question: Where did fortune cookies origin ate ? Answer : Fortune cookies originated in San Francisco

next token:6662- San; token_probit: 0.9997644221069274; running_probit: 1.9627333831057165e-18

Question: Where did fortune cookies origin ate ? Answer : Fortune cookies originated in San

'San Francisco' Running Probit ~= 2e-18
'United States' Running Probit ~= 2e-13

It knows the fact, but the earlier generative styling starts moving it towards the Invention/United States response first
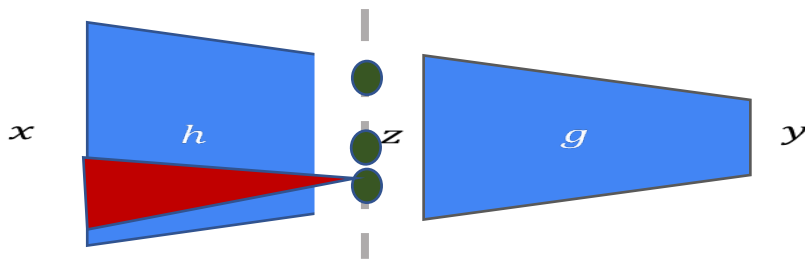
42

# Internal Influence

# Influence-directed explanations [Leino, Sen, Datta, Fredrikson, Li 2018]

Explaining property of a ML system =
**identify influential factors +
make them human interpretable**

- Influence: What are important factors causing this model property?
- Interpretation: What do these factors mean?

# Influence-directed explanations for deep networks

- Rank causally influential neurons in internal layers (novel!)
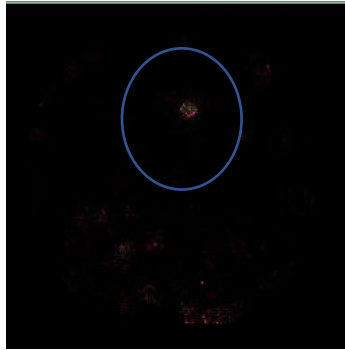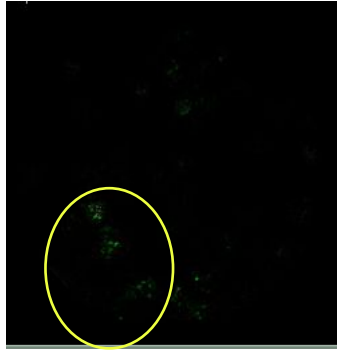- Give them interpretation using visualization techniques (prior work)
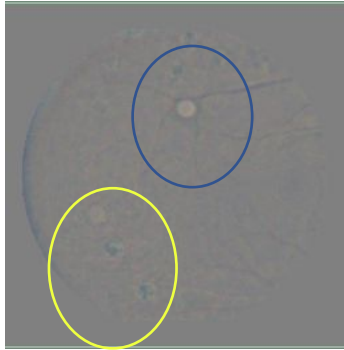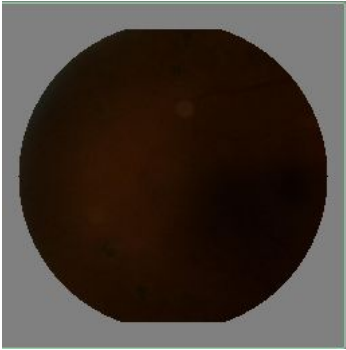


First result with internal influence measure for deep networks

# Why classified as diabetic retinopathy stage 5?
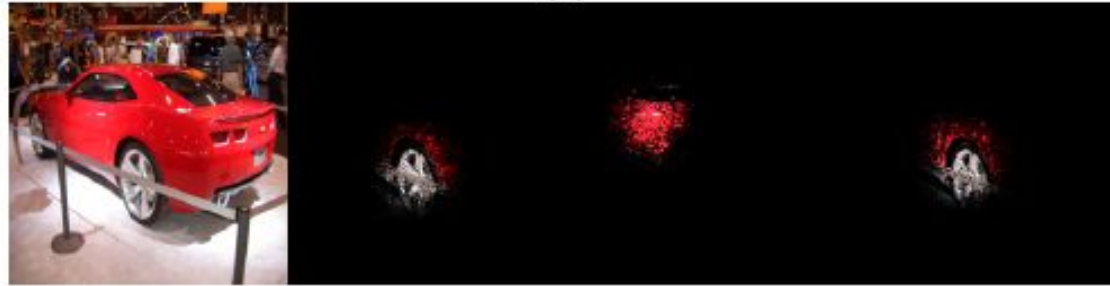
Inception network

Optic disk

Lesions

# Why did the network classify input as sports car?



Input image                    Influence-directed Explanation

# Why sports car instead of convertible?

VGG16 ImageNet model



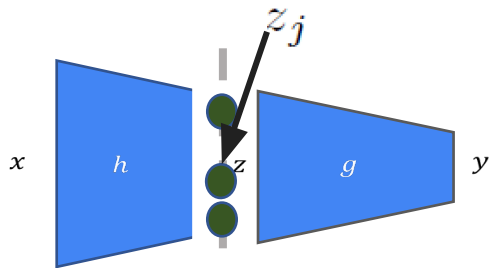Input image                      Influence-directed Explanation

Uncovers high-level concepts that generalize across input instances

# Distributional influence

Influence = average gradient over distribution of interest



$$y = f(x) = g(h(x))$$

$$\chi_j^s(f, P) = \int_{\mathcal{X}} \left.\frac{\partial g}{\partial z_j}\right|_{h(\mathbf{x})} P(\mathbf{x}) d\mathbf{x}$$

Gradient

Weighted by probability of input x

49

For input x  [note z = h(x)]

Theorem: Unique measure that satisfies a set of natural properties

# Interpreting influential neurons



Depicts interpretation (visualization) of 3 most influential neurons

- Slice of VGG16 network: conv4_1
- Inputs drawn from distribution of interest: delta distribution
- Quantity of interest: class score for correct class

# Interpreting influential neurons



Visualization method: Saliency maps [Simonyan et al. ICLR 2014]

- Compute gradient of neuron activation wrt input pixels
- Scale pixels of original image accordingly

# Internal Explanations via Influence Paths



- Influence paths provide insights into misclassifications
- Model can be compressed down the influential paths without changing the utility of the model

**Influence Paths**
Lu, Mardziel, Leino, Fedrikson, Datta, ACL '20

# Internal Explanations via Influence Patterns
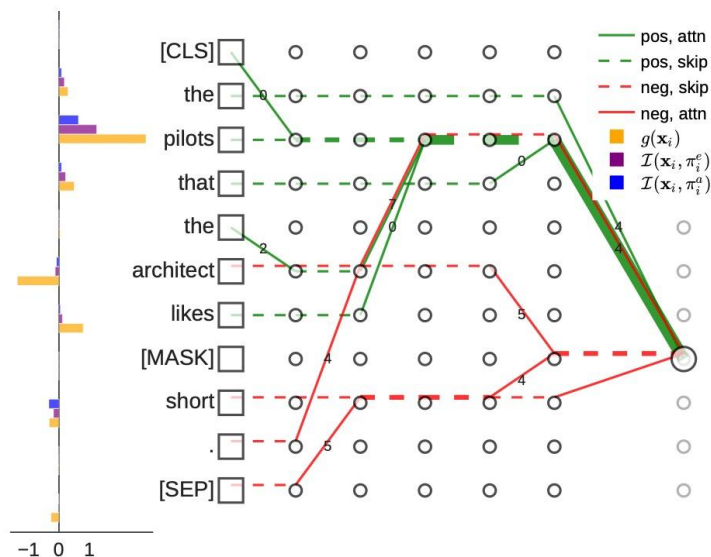
## BERT v.s. LSTM

- Scaling up method to identify influential paths

- Prevalence of "copy" and "transfer" operations to carry context



**Influence Patterns for BERT**
Lu, Wang, Mardziel, Datta, NeurIPS 2021

# What-If Exploration

# What-If Exploration

**Probe the model on various What-If scenarios.**

- Examples:
    - What if "he" was replaced with "she"
    - What if we add a punctuation at the end of the sentence
- **Intuitive:** What you see is what you get
- **Highly expressive:** Most explainability techniques are a summarization of what-if behavior

**Applications:**

- Model understanding / debugging
- Algorithmic Recourse
- Prompt design

# What-If Exploration

**But,**

- How do we navigate the (vast) space of what-if scenarios?

- How do we identify what-if scenarios that achieve a target prediction?

# Problem Statement

**Given an input and a prediction target, identify a set of minimal perturbations that achieve the target**

- Perturbations defined replacing features with empty value (e.g., drop a word) or replacing them with a reference feature

- Minimality is defined via **partial order (≼)** on the space of perturbations

    - E.g., perturbation {he → she} is more preferable (≼) to {he → she, him → her}

# Technique: Targeted What-Ifs

- Iterate through the space of perturbation in topologically sorted order

- Return perturbations that achieve the prediction target with at least probability $\tau$

# Technique: Targeted What-Ifs

- Iterate through the space of perturbation in topologically sorted order

- Return perturbations that achieve the prediction target with at least probability $\tau$

**Paper:** [Local Explanations via Necessity and Sufficiency: Unifying Theory and Practice](#), UAI 2021

- Frames the problem using the theory of sufficient and necessary causes, and proves a correctness guarantee

  - **[Soundness]** All returned perturbations are minimal and achieve the target

  - **[Completeness]** All target-achieving, minimal perturbations are returned

# Case study from a Search team: Detecting Irrelevant Features

**Issue:** A search model was predicting high pCTR for certain queries paired with an irrelevant result.

**Debugging:** Identify query token ablations (what-ifs) that lowered the pCTR

**Finding:** Perturbations identified out-of-vocab (OOV) tokens, e.g., the token "ph8" in query "water filter ph8"

**Root cause:** Model was not trained well on queries with OOV tokens.

**Fix:** Increase the vocab frequency threshold (so that more OOV tokens are seen during training) and retrain. This fixed the issue!

# LLM Application Idea: Prompt Perturbations

**Prompt:** I am going to show you a query and top-3 search results for the query. Please provide a concise answer to the query based on the search results. Do not use any information outside the search results. The answer must be no longer than 3 sentence. You may return "irrelevant results" if the search results do not contain an answer to the question

Query: <Query>

Search Results: <Search Results>

**Is the model taking all instructions into account?**

**Is the model fixating on unimportant aspects of the prompt like spaces and punctuation?**

# LLM Application Idea: Prompt Perturbations

Prompt: `I am going to show you a query and top-3 search results for the query. Please provide a concise answer to the query based on the search results. Do not use any information outside the search results. The answer must be no longer than 3 sentence. You may return "irrelevant results" if the search results do not contain an answer to the question`

`Query: <Query>`

`Search Results: <Search Results>`

**Is the model taking all instructions into account?**

**Is the model fixating on unimportant aspects of the prompt like spaces and punctuation?**
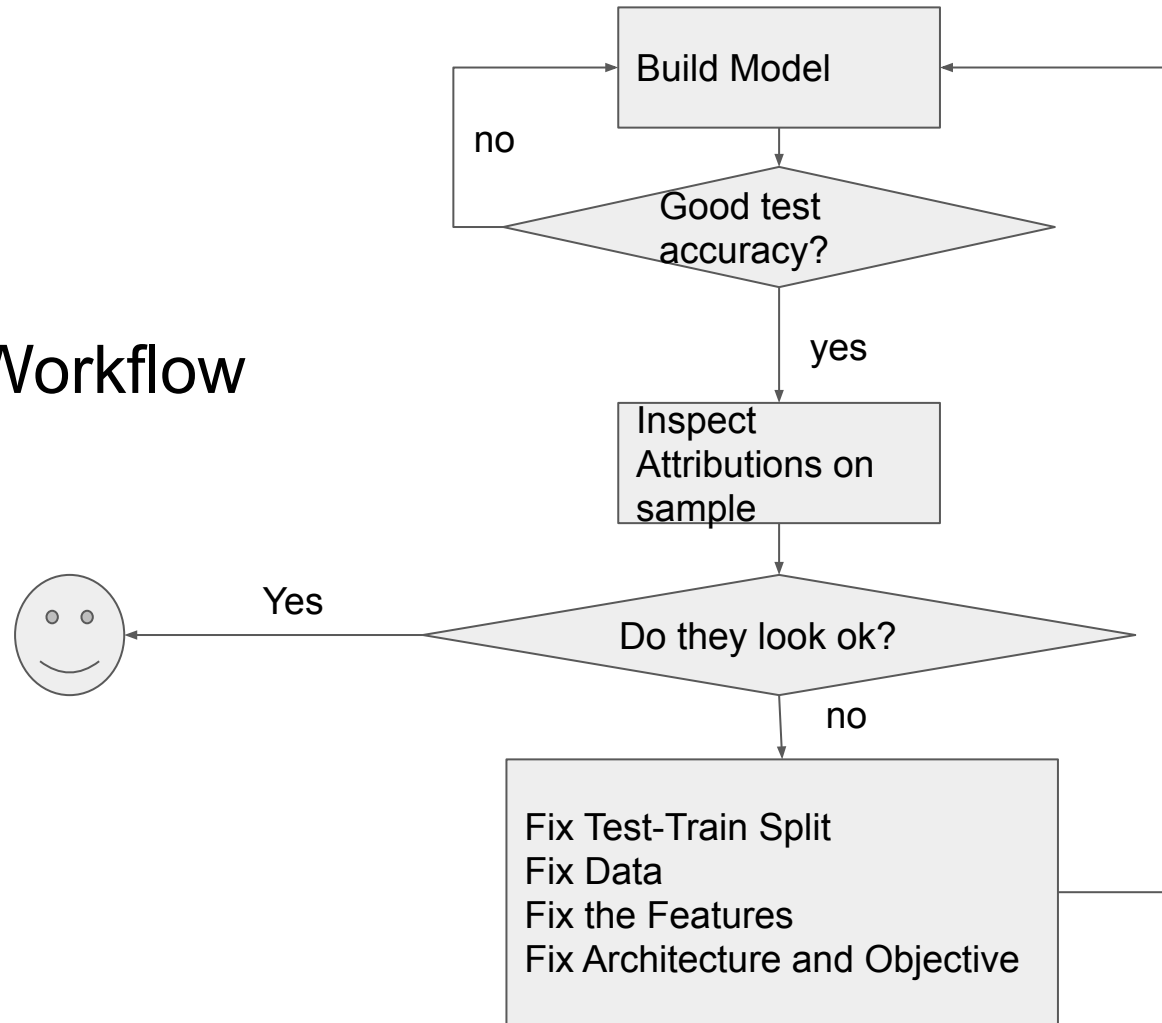
**Applying the method:**

- Define a (ordered) space of perturbations

- Identify maximal token perturbations that preserve its performance (on an input set)

- Identify minimal token perturbations that drastically alter its performance (on an input set)
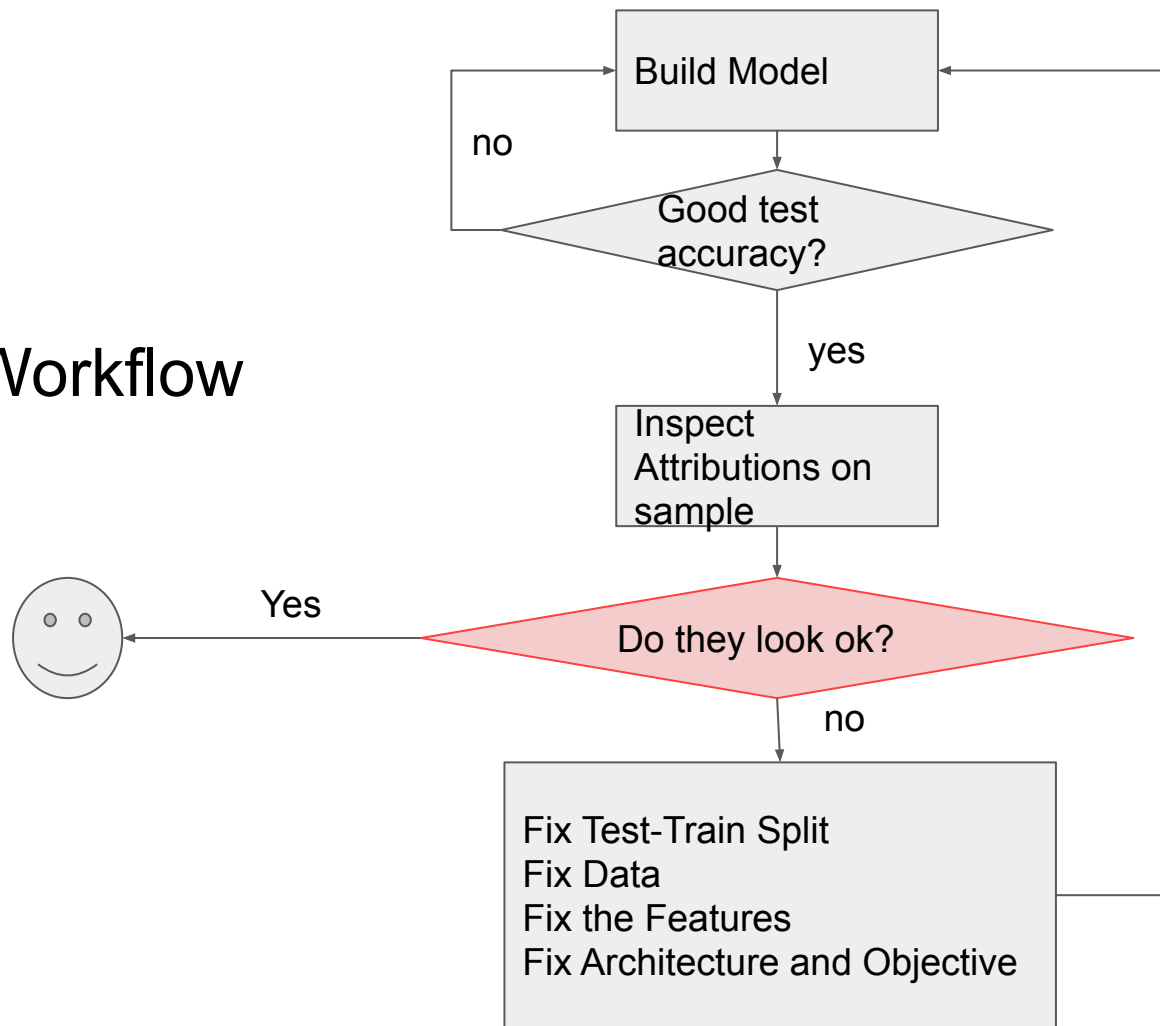
# Summary

- Test accuracy alone can be misleading
  - Examine model performance on slices
  - Assess if test set is representative of deployment
- Probe the model's reasoning on individual predictions
  - Is the model relying on spurious/irrelevant features?
  - Is the model ignoring relevant features?
- Feature attributions are surprisingly good at uncovering model behaviors
  - Proper visualization + Human thought is crucial for turning attributions to insights.

Some limitations and caveats for feature attributions

Debugging Workflow

Build Model

Good test accuracy?

no

yes

Inspect Attributions on sample

Do they look ok?

Yes

no

Fix Test-Train Split
Fix Data
Fix the Features
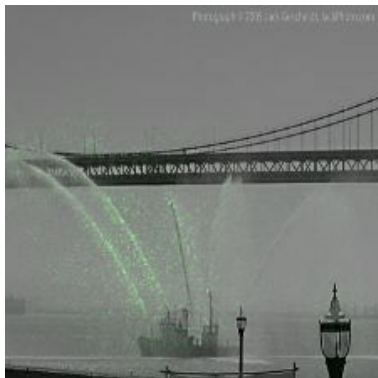Fix Architecture and Objective

Debugging Workflow
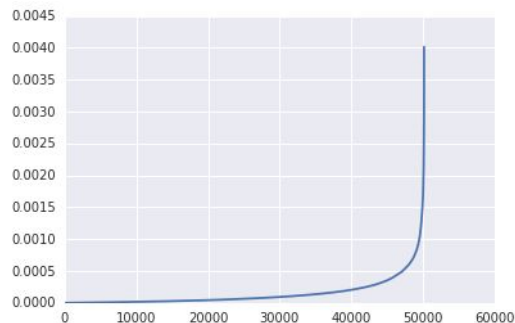
# Role of the Human Analyst

- Humans are poor at foreseeing problems

- Humans excel at understanding real world implications of specific explanations

  - Disease prediction: "Pen marks won't be available on X-rays in deployment"

  - Question answering: "most words in a question matter"

- Proper visualization is very important in making attributions intelligible to humans

# Importance of Visualization

**Naive** scaling of attributions from 0 to 255

Attributions have a **large range** and **long tail** across pixels

**After clipping** attributions at 99% to reduce range



**Paper:** Exploring Principled Visualizations for Deep Network Attributions, IUI Workshop 2019

# Feature Attributions are pretty shallow

Attributions do not explain:

- How the network combines the features to produce the answer?

- What training data influenced the prediction

- Why gradient descent converged

- etc.

**Attributions are useful when the network behavior entails that a strict subset of input features are important**

# Evaluating Integrated Gradients
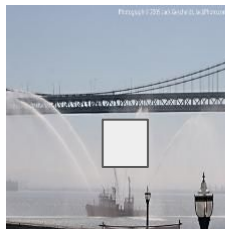
# Evaluating Integrated Gradients

- Ablate top attributed features and examine the change in prediction
  - <u>Issue</u>: May introduce artifacts in the input (e.g., the square below)



- Compare attributions to (human provided) groundtruth on "feature importance"
  - <u>Issue 1</u>: Attributions may appear incorrect because the network reasons differently
  - <u>Issue 2 </u>: **Confirmation bias**

# Evaluating Integrated Gradients

- Ablate top attributed features and examine the change in prediction
  - <u>Issue</u>: May introduce artifacts in the input (e.g., the square below)



- Compare attributions to (human provided) groundtruth on "feature importance"
  - <u>Issue 1</u>: Attributions may appear incorrect because the network reasons differently
  - <u>Issue 2</u> : **Confirmation bias**

The mandate for attributions is to be faithful to the model's reasoning

# Our Approach: Axiomatic Justification

- List **desirable criteria (axioms)** for an attribution method

- Establish a uniqueness result: X is the **only** method that satisfies these criteria

# Axioms

- **Insensitivity:** A variable that has no effect on the output gets no attribution

- **Sensitivity**: If baseline and input differ in a single variable, and have different outputs, then that variable should receive some attribution

- **Linearity preservation**: Attributions(α*F1 + ß*F2) = α*Attributions(F1) + ß*Attributions(F2)

- **Implementation invariance**: Two networks that compute identical functions for all inputs get identical attributions

- **Completeness**: Sum(attributions) = F(input) - F(baseline)

- **Symmetry**: Symmetric variables with identical values get equal attributions

# Result

**Theorem** [ICML 2017]: Integrated Gradients is the **unique** path-integral method satisfying: Sensitivity, Insensitivity, Linearity preservation, Implementation invariance, Completeness, and Symmetry

**Historical note:**

- Integrated Gradients is the **Aumann-Shapley method** from cooperative game theory, which has a similar characterization; see [Friedman 2004]