

# CSCE 771: Computer Processing of Natural Language

## Lecture 24: Ethical Concerns with NLP, Trusted AI and Societal Impact

---

PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

15<sup>TH</sup> NOVEMBER, 2022

***Carolinian Creed: “I will practice personal and academic integrity.”***

# Organization of Lecture 24

---

- Opening Segment
  - Announcements

- Main Lecture



- Concluding Segment
  - About Next Lecture – Lecture 25

## Main Section

- The issue of trust
- Ethics and fairness issues
- Mitigation
  - Explanation and interpretation
  - Rating for trust

# Recent Classes

---

Nov 1 (Tu)	NLP Task: Sentiment
Nov 3 (Th)	NLP Task: Summarization
Nov 8 (Tu)	
Nov 10 (Th)	Conversation Agents
Nov 15 (Tu)	Ethical Concerns with NLP, Trusted AI and Societal Impact
Nov 17 (Th)	Working with LLMs for NLP Tasks - programming, Quiz
Nov 22 (Tu)	Paper presentations
Thanksgiving Holiday	
Nov 29 (Tu)	Project presentations
Dec 1 (Th)	Project presentations
Dec 8 (Tu)	Quiz

## Review of Lecture 22

- Different types of chatbots
- Potential for using them
- Different ways of building them
  - Rule based methods
  - (Deep) learning based methods
- Applications
- Ethical Issues

# Announcements

---

# Reference: Project Rubric

---

- **Project results – 60%**
  - Working system ? – 30%
  - Evaluation with results superior to baseline? – 20%
  - Considered related work? – 10%
- **Project efforts – 40%**
  - Project report – 20%
  - Project presentation (updates, final) – 20%
- **Bonus**
  - Challenge level of problem – 10%
  - Instructor discretion – 10%
- **Penalty**
  - Lack of timeliness as per announced policy (right) - up to 60%

## Milestones

- Penalty: **not** ready by Sep 15, 2022 **[-20%]**
- Project report **not** ready by Nov 10, 2022 **[-20%]**
- Project presentations not ready by Nov 15, 2022 **[-10%]**

Project presentation DUE today!

# Deadlines for Project Reports and Presentations

---

- Since the deadlines were posted since the beginning of the semester, we will not move them. However, submissions made until respective deadline can be updated till Nov 20.
- For Reports:
  - Prepare an initial version of the report by deadline (Nov 10, 2022), put in your GitHub, send me a note. It should be complete in terms of all sections and initial content.
  - You can update and post new copies until Nov 20. If updating, make sure to not overwrite the initial versions. If I do not see the initial version, I will have to penalize for missed deadline.
- For Presentations:
  - Prepare an initial version of the report by deadline (Nov 15, 2022), put in your GitHub, send me a note. It should be complete in terms of all sections and initial content.
  - You can update and post new copies until Nov 20. If updating, make sure to not overwrite the initial versions. If I do not see the initial version, I will have to penalize for missed deadline.

# Project Report Guidelines

---

- Use template of ACM Computing Surveys – Latex or Word - <https://www.acm.org/publications/authors/submissions>
- Consider your report as a paper. Sections to have will be similar
  - **Abstract**: 1-line each on what, how, result // Optional
  - **Introduction**: motivation for the work // Optional
  - **Problem** // Clearly state input and output
  - **Related Work** // What are closely related work?
  - **Approach** // How does your system work?
  - **Evaluation** // How is the result better than a baseline? What better could have been done ?
  - **Discussion** // About results, what more could be done, anything else interesting
  - **Conclusion** // Optional
  - **References**

# Project Presentation Template (2 mins)

---

- Project Name
- Problem
- Approach
- Result
- Comment:
  - Challenges faced
  - Need help

- Test Case – *how your program run*
- Evaluation



# Other Milestones / Deadlines

---

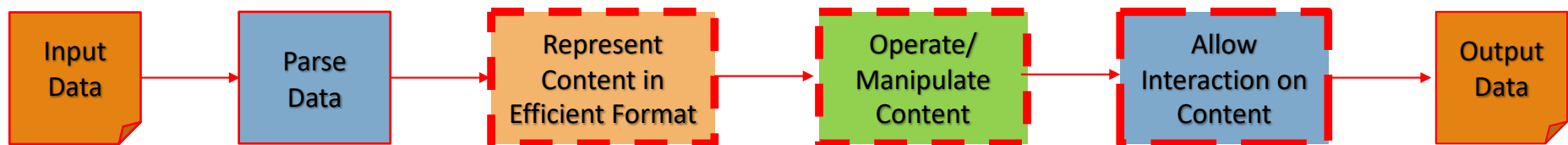
- Nov 17: Quiz 3
  - Programming quiz; Use class time to review material
- Nov 22: Research paper presentation (2 mins)
  - Add paper title and venue (column M) – has to be a research or application paper in last two years (2020-2022) at a top AI/NLP/Image/Audio conference: ACL, AAAI, IJCAI, NeurIPS, CIKM, CVPR, ICML, WWW.
  - Prepare 1-slide summary containing the following and present in 2 mins in class
    - **Summary:** problem, solution, related work, evaluation, contributions
    - **Opinion:** What you liked or did not like

# Main Lecture

---

# Language Processing – Remain Trustworthy

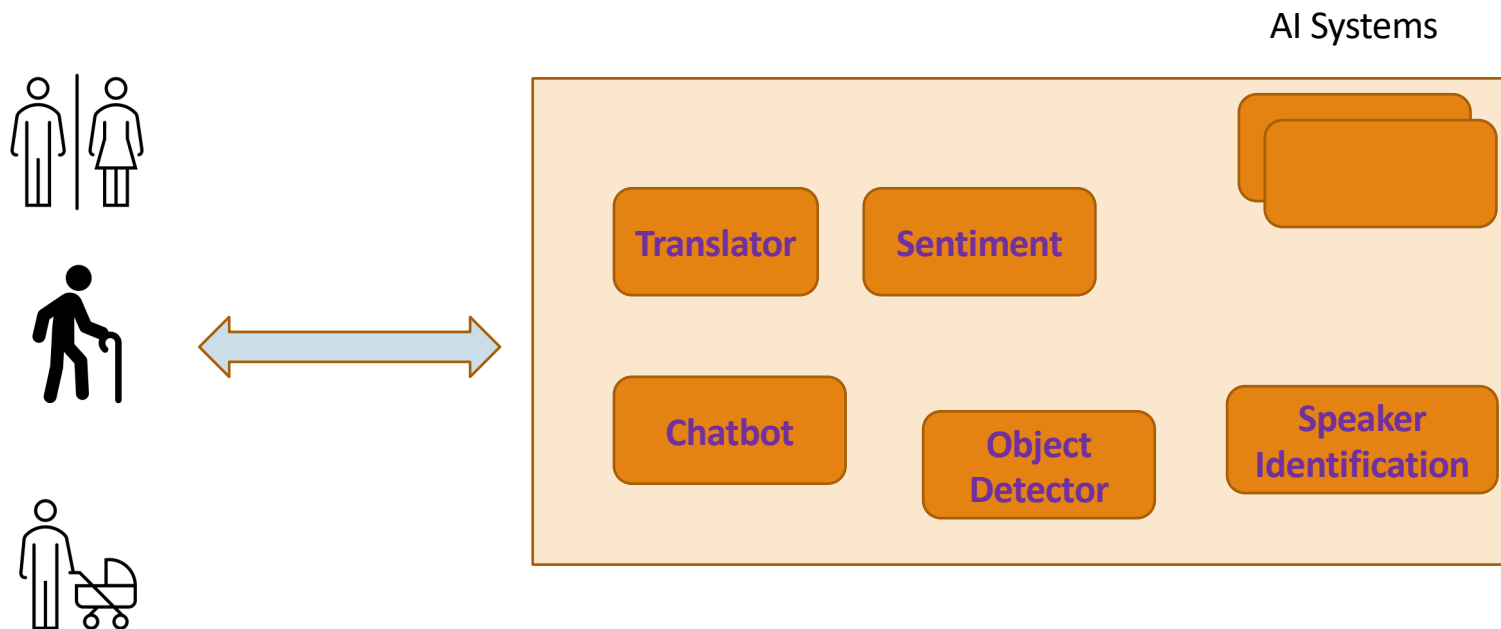
---



# The Problem of Trust

---

# Technology and People



**Trust:** *Can people trust AI systems to perform capably, consistently, and with human values?*

# What are the Components of Trust (Technology)

---

1. Competent – does what it is supposed to do
2. Reliable – including, well tested
3. Upholds human values, social good
  1. Fairly and ethically used
  2. Adequate data management & preserves privacy
4. Allows human-technology interaction
  1. Explainable, transparent
  2. How does the system give its result?

**Reference:** Trustworthy Machine Learning, Kush R. Varshney, 2022  
<http://www.trustworthymachinelearning.com/>

# Components of Trust for AI

1. Competent – does what it is supposed to do
2. Reliable – including, well tested
3. Upholds human values
  1. Fairly and ethically used
  2. Adequate data management & preserves privacy
4. Allows human-technology interaction
  1. Explainable, transparent
  2. How does the system give its result?

	AI – Word Tag Cloud	AI – Image Search	AI – Self-driving Car	AI-powered Chatbot: Medical Guide
Competent	x	x	?	x
Reliable	x	?	?	?
Upholds human values	?	?	?	?
Allows human interaction	x	x	?	?

x: yes; -: not applicable; ?: questionable

## Illustration: A Seemingly Innocuous Chatbot

### Potential Issues

- Leak information
- Abusive language
- Complex response

#### References:

1. Ramashish Gaurav, Biplav Srivastava, Estimating Train Delays in a Large Rail Network Using a Zero Shot Markov Model, IEEE International Conference on Intelligent Transportation Systems (ITSC). On Arxiv at: <https://arxiv.org/abs/1806.02825>, June 2018 [Train delay, prediction]
2. Himadri Mishra, Ramashish Gaurav, Biplav Srivastava, Train Status Assistant for Indian Railways, On Arxiv at: <https://arxiv.org/abs/1809.08509>, Sep 2018, Video: <https://www.youtube.com/watch?v=a-ABv29H6XU> [Chatbot, Train delay assistant]

## TDEBot



TDEBot, 3:29 PM

Train Number 12312 will be delayed by 278.0 minutes at HWH station on 2018-10-18



TDEBot, 3:29 PM

The bottleneck station is FTP causing delay of 90.2 minutes on 2018-10-18



TDEBot, 3:32 PM

Sorry, I didn't understand! Please Try again



TDEBot, 3:33 PM

Train 12312 will not be mitigated any more after station ALD on 2018-10-18. It will arrive even later by 52.0 minutes

is train 12312 on time today?

3:29 PM

Where is the bottleneck?

3:32 PM

What is FTP?

3:32 PM

What is the delay at Allahabad?



# UK's National Screening Committee Assessment on Use of AI for Breast Screening

*"The current review looked at the evidence on:*

Details: <https://t.co/6RAgE5eBCH>, Feb 2022

- how good AI is at finding cancers in breast cancer screening*
- what benefits and harms AI has for the women who are screened or for the screening program and the health professionals involved*

*Based on the current evidence, the **UK NSC does not recommend using AI in the NHS breast cancer screening program**. This is because:*

- the use of AI systems would change the current screening program\* therefore it is important to assess how accurate AI is in breast screening clinical practice before changing it*
- the performance of AI systems varies in different settings but there are no good quality studies in the UK*
- it is unclear how good AI is at finding different types of breast cancer or at finding breast cancers in different groups of women (for example different ethnic groups)*
- AI might reduce the workload of staff, the number of cancers missed at screening, and the number of women called back for further tests when they do not have cancer, however, the quality of evidence is very low. "*

\* Changed spelling

# Instability of AI is Well Recorded

---

[Text] [Su Lin Blodgett](#), [Solon Barocas](#), [Hal Daumé III](#), [Hanna Wallach](#), Language (Technology) is Power: A Critical Survey of “Bias” in NLP, Arxiv - <https://arxiv.org/abs/2005.14050>, 2020 [NLP Bias]

[Image] Vegard Antun, Francesco Renna, Clarice Poon, Ben Adcock, and Anders C. Hansen, On instabilities of deep learning in image reconstruction and the potential costs of AI, <https://doi.org/10.1073/pnas.1907377117>, PNAS, 2020

[Audio] Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R. Rickford, Dan Jurafsky, and Sharad Goel, Racial disparities in automated speech recognition, PNAS April 7, 2020 117 (14) 7684-7689, <https://doi.org/10.1073/pnas.1915768117>, March 23, 2020

# Current AI: Capabilities, Limitations, Ethical issues

## Capabilities

### Machine Learning

- Learning from data (Deep, Reinforced, Supervised/Unsupervised/Self Supervised)
- Hidden patterns in huge amounts of data
  - Prediction, perception tasks
  - Correlation, pattern discovery, data mining
- Flexible, can handle uncertainty

### Rule-based, symbolic, and logical approaches

- Explicit procedure to solve a problem
- Reasoning, planning, scheduling, optimization for complex problems
- Symbolic, traceable, explainable

### Limitations

- Generalizability and Abstraction
- Robustness and Resiliency
- Contextual awareness
- Multi-agent cooperation
- Resource efficiency (examples, energy, computing power)
- Adaptability
- Causality

### AI ethics issues

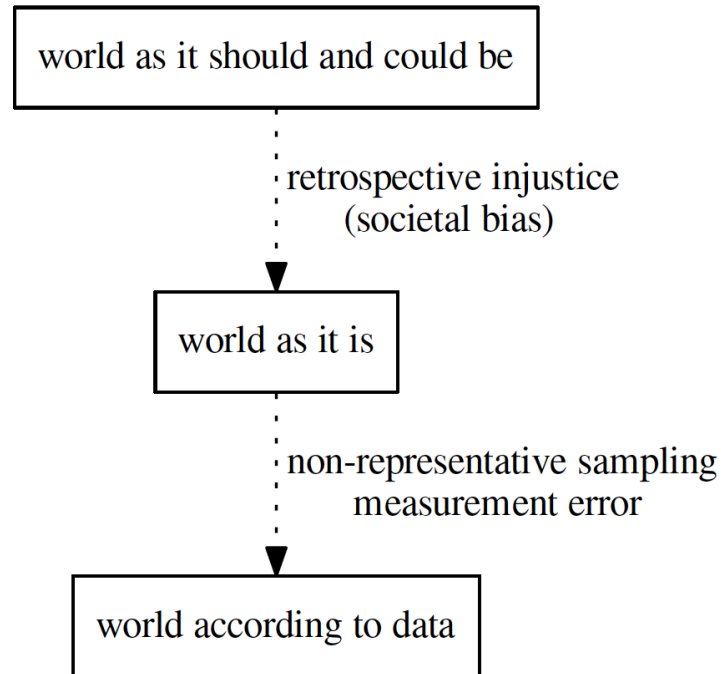
- Trust
  - Fairness, robustness, explainability, causality, transparency
- Data governance, privacy, liability, human agency, impact on work and society
- AI autonomy vs augmented intelligence
- Real vs online life, metrics of success/goals

**Slide credit:** Francesca Rossi

# Ethics and Fairness

---

## Data, Bias and the World We Live In



from "Prediction-Based Decisions and Fairness" by Mitchell, Potash and Barocas, 2018

**when data is about people, bias can lead to discrimination**

## Usage of People-Neutral Technology for People-Sensitive Applications

**Data science is algorithmic, therefore it cannot be biased!** And yet...

- All traditional evils of **discrimination**, and many new ones, exhibit themselves in the data science eco system
- **Bias** that is inherent in the data or in the process, and that is often due to systemic discrimination, is propelled and amplified
- **Transparency** helps prevent discrimination, enable public debate, establish **trust**
- Technology alone won't do: also need **policy**, **user involvement** and **education**



<http://www.allenoverly.com/publications/en-gb/Pages/Protected-characteristics-and-the-perception-reality-gap.aspx>

# What is Specific to AI?

- AI needs **data**
  - Data privacy and governance
- AI is often a **black box**
  - Explainability and transparency
- AI can make **decisions/recommendations**
  - Fairness and value alignment
- AI is based on statistics and has always a small percentage of **error**
  - Who is accountable if mistakes happen?
- AI can infer our preferences and **manipulate** them
  - Human and moral agency
- AI is very **pervasive and dynamic**
  - Larger negative impacts for tech misuse
  - Fast transformation of jobs and society

## Credits:

Tutorial on [Trusting AI by Testing and Rating Third Party Offerings at IJCAI 2020](#), Biplav Srivastava, Francesca Rossi, Jan 2021

# AI Ethics

---



Multidisciplinary field of study



How to optimize AI's beneficial impact while reducing risks and adverse outcomes



How to design and build AI systems that are aware of the values and principles to be followed in the deployment scenarios



To identify, study, and propose technical and nontechnical solutions for ethics issues arising from the pervasive use of AI in life and society

**Credits:**

Tutorial on [Trusting AI by Testing and Rating Third Party Offerings](#) at IJCAI 2020, Biplav Srivastava, Francesca Rossi, Jan 2021



# A Tale of Two Definitions

---

## Machine Learning

- Often refers to members of protected classes as those in “minority and marginalized groups”
- Analysis of demographics data can lead to better anti-discrimination policies

## Legal

- Focus on equal treatment, regardless of attributes such as race and gender
- Landmark affirmative action cases have concluded that schools seeking to increase racial diversity cannot use racial quotas or point systems.

**Source:** To Prevent Algorithmic Bias, Legal and Technical Definitions around Algorithmic Fairness Must Align,  
<https://www.partnershiponai.org/to-prevent-algorithmic-bias-legal-and-technical-definitions-around-algorithmic-fairness-must-align/>  
**Paper:** <https://arxiv.org/pdf/1912.00761.pdf>

# Examples of Computational / AI Services and Bias

---

Search results, e.g., matching (jobs), nearest (hospitals, taxi-ride, groceries)

- **Some possible biases:** age, gender, racial, income
- **Impact :** failure to be diverse in employment (match), deny or provide costlier services where most needed

Language translator

- **Some possible biases:** gender, religious, racial
- **Impact:** failure to recognize gender may lead to selection of wrong/indecent phrase in target language which can cause uproar

Medical condition detector

- **Some possible biases:** gender, racial
- **Impact :** failure to recognize entities may lead to mis-diagnosis

Image caption generator

- **Some possible biases:** Sexual, religious, racial
- **Impact :** failure to recognize entities in image may lead to selection of wrong phrases and generation of wrong/indecent caption which can cause uproar

# Main AI Ethics Issues



DATA GOVERNANCE  
AND PRIVACY



FAIRNESS AND  
INCLUSION



HUMAN AND  
MORAL AGENCY



VALUE ALIGNMENT



ACCOUNTABILITY



TRANSPARENCY AND  
EXPLAINABILITY



TECHNOLOGY  
MISUSE

**Credits:**

Tutorial on [Trusting AI by Testing and Rating Third Party Offerings](#) at IJCAI 2020, Biplav Srivastava, Francesca Rossi, Jan 2021

# Does Trust Matter – A Recent IBM IBV Survey

1,250 global executives in late 2018:  
Representing 20 industries and over 26 countries on 6 continents, including members of boards of directors, chief executive officers (CEOs), chief information officers (CIOs), chief technology officers (CTOs), chief data officers (CDOs), chief human resource officers (CHROs), chief risk officers (CROs), general counsels, and government policy officials.

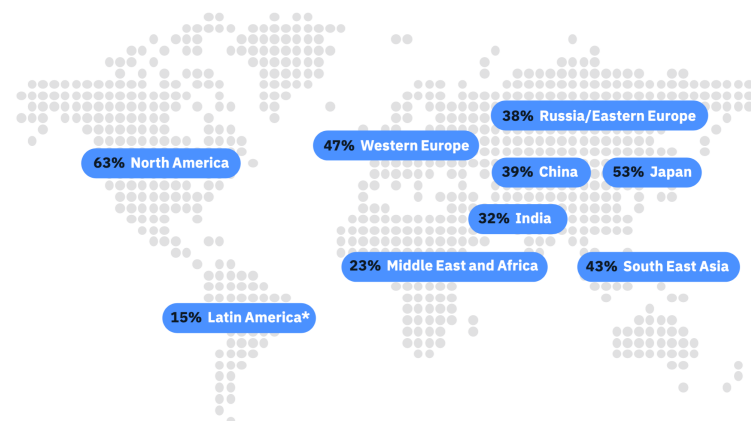
## Questions:

- Who is responsible for helping ensure that ethics are integrated into AI, within corporations and outside?
- What is most important in ethically harnessing the power of AI?
- And how can society best use AI for good

## Insights

- 81% of consumers say they became more concerned over the prior year with how companies use their data, and 75% percent are now less likely to trust organizations with their personal information
- Well over half of all executives point to the CTO and CIO as primarily accountable for AI ethics.
- Executives expect technology firms will greatly influence AI ethics, followed by governments and customers – with other companies last on the list.
- Three main areas of ethical risk: data responsibility, value alignment, and algorithmic accountability

The importance organizations place on AI ethics varies across regions



Advancing AI ethics beyond compliance From principles to practice,  
April 2020; PDF: <https://www.ibm.com/downloads/cas/J2LAYLOZ>

\*Count is less than 20.  
Source: 2018 IBM Institute for Business Value Global AI Ethics Study. Q: Importance of AI ethics in your organization, N=1,247.

Source: *Fairness and Machine Learning* by Solon Barocas, Moritz Hardt, Arvind Narayanan (<https://www.fairmlbook.org>)

# A Step Towards Fairness

## Broad classes

- **Individual fairness:** similar individuals to be treated similarly
- **Group fairness:** statistical property of decision as a group should be representative of the population
- **Both individual and group fairness, and use a single metric:** generalized entropy index

**Guidance:** *Selection of metric is application driven*

Name	Closest relative	Note	Reference
Statistical parity	Independence	Equivalent	Dwork et al. (2011)
Group fairness	Independence	Equivalent	
Demographic parity	Independence	Equivalent	
Conditional statistical parity	Independence	Relaxation	Corbett-Davies et al. (2017)
Darlington criterion (4)	Independence	Equivalent	Darlington (1971)
Equal opportunity	Separation	Relaxation	Hardt, Price, Srebro (2016)
Equalized odds	Separation	Equivalent	Hardt, Price, Srebro (2016)
Conditional procedure accuracy	Separation	Equivalent	Berk et al. (2017)
Avoiding disparate mistreatment	Separation	Equivalent	Zafar et al. (2017)
Balance for the negative class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Balance for the positive class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Predictive equality	Separation	Relaxation	Chouldechova (2016)
Equalized correlations	Separation	Relaxation	Woodworth (2017)
Darlington criterion (3)	Separation	Relaxation	Darlington (1971)
Cleary model	Sufficiency	Equivalent	Cleary (1966)
Conditional use accuracy	Sufficiency	Equivalent	Berk et al. (2017)
Predictive parity	Sufficiency	Relaxation	Chouldechova (2016)
Calibration within groups	Sufficiency	Equivalent	Chouldechova (2016)
Darlington criterion (1), (2)	Sufficiency	Relaxation	Darlington (1971)

# German Credit Data

- Dataset that classifies people's credit risk based on their individual attributes such as Age, Income, Gender, etc.

Example Instance:

A11 6 A34 A43 1169 A65 A75 4 A93 A101 4 A121 67 A143 A152 2 A173 1 A192 A201 1

<https://archive.ics.uci.edu/ml/datasets/Statlog+%28German+Credit+Data%29>

- Each entry represents an individual who takes credit from a bank
- Each entry is classified as *Good* or *Bad* credit risk based on their profile
  - It is worse to classify a **customer as good when they are bad**, than it is to classify a **customer as bad when they are good**.
- 1000 rows of data, each with 20 attributes to check bias against

Dua, D. and Graff, C. (2019). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>], Irvine, CA: University of California, School of Information and Computer Science

# Picking the Appropriate Fairness Metrics

**Statistical Parity Difference:** Difference of the rate of favorable outcomes received by the unprivileged group to the privileged group

**Equal Opportunity:** Difference of true positive rates between the two groups

**Average Odds Difference:** Difference of false positive rate and true positive rate between the groups

**Disparate Impact:** The ratio of rate of favorable outcome for the unprivileged group to that of the privileged group

**Theil Index:** The generalized entropy of benefit for all individuals in the dataset, with  $\alpha = 1$

Name	Closest relative	Note	Reference
Statistical parity	Independence	Equivalent	Dwork et al. (2011)
Group fairness	Independence	Equivalent	
Demographic parity	Independence	Equivalent	
Conditional statistical parity	Independence	Relaxation	Corbett-Davies et al. (2017)
Darlington criterion (4)	Independence	Equivalent	Darlington (1971)
Equal opportunity	Separation	Relaxation	Hardt, Price, Srebro (2016)
Equalized odds	Separation	Equivalent	Hardt, Price, Srebro (2016)
Conditional procedure accuracy	Separation	Equivalent	Berk et al. (2017)
Avoiding disparate mistreatment	Separation	Equivalent	Zafar et al. (2017)
Balance for the negative class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Balance for the positive class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Predictive equality	Separation	Relaxation	Chouldechova (2016)
Equalized correlations	Separation	Relaxation	Woodworth (2017)
Darlington criterion (3)	Separation	Relaxation	Darlington (1971)
Cleary model	Sufficiency	Equivalent	Cleary (1966)
Conditional use accuracy	Sufficiency	Equivalent	Berk et al. (2017)
Predictive parity	Sufficiency	Relaxation	Chouldechova (2016)
Calibration within groups	Sufficiency	Equivalent	Chouldechova (2016)
Darlington criterion (1), (2)	Sufficiency	Relaxation	Darlington (1971)

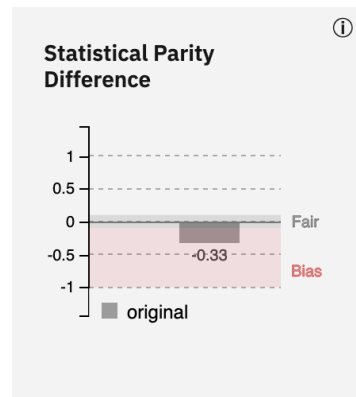
# Checking Bias Metrics: *Age*

## Protected Attribute: Age

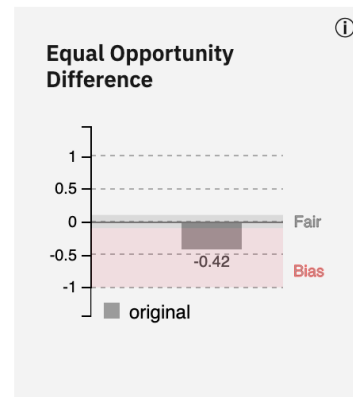
Privileged Group: **Old**, Unprivileged Group: **Young**

Accuracy with no mitigation applied is 76%

With default thresholds, bias against unprivileged group detected in 4 out of 5 metrics

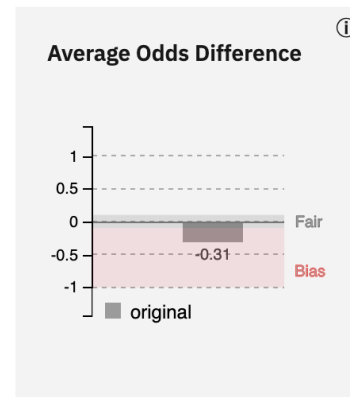


Ideal Value: 0  
Fairness Range: (-0.10, 0.10)



Ideal Value: 0  
Fairness Range: (-0.10, 0.10)

A value of < 0 implies higher benefit for the privileged group and a value > 0 implies higher benefit for the unprivileged group.



Ideal Value: 0  
Fairness Range: (-0.10, 0.10)



Ideal Value: 1.0  
Fairness Range: (0.80, 1.20)

A value < 1 implies higher benefit for the privileged group and a value > 1 implies a higher benefit for the unprivileged group.



Ideal Value: 0  
Fairness is indicated by lower scores, higher scores represent inequality



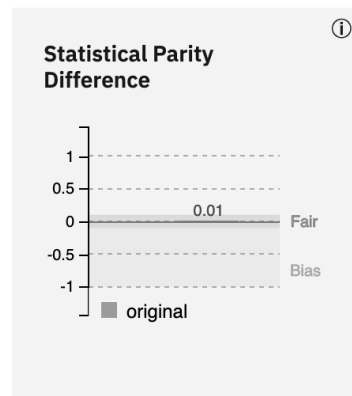
# Checking Bias Metrics: *Gender*

## Protected Attribute: Sex

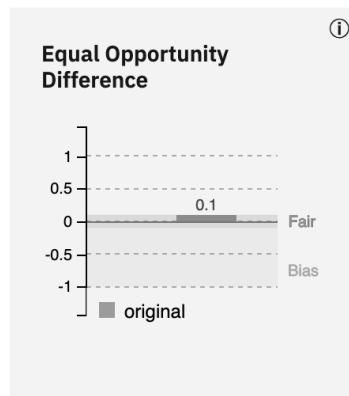
Privileged Group: **Male**, Unprivileged Group: **Female**

Accuracy with no mitigation applied is 76%

With default thresholds, bias against unprivileged group detected in 0 out of 5 metrics

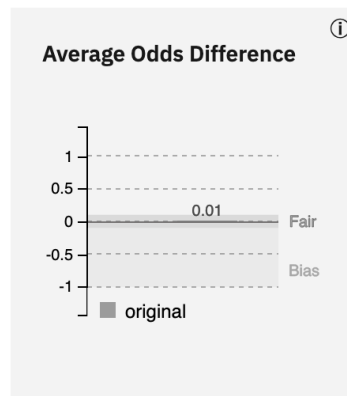


Ideal Value: 0  
Fairness Range: (-0.10, 0.10)

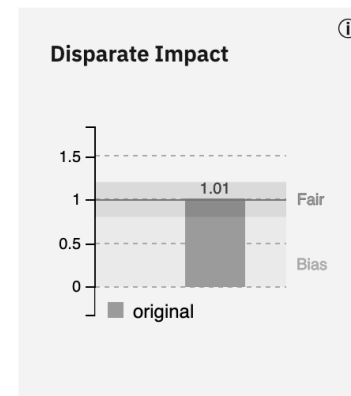


Ideal Value: 0  
Fairness Range: (-0.10, 0.10)

A value of  $< 0$  implies higher benefit for the privileged group and a value  $> 0$  implies higher benefit for the unprivileged group.

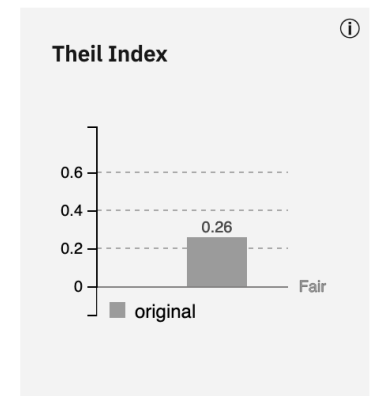


Ideal Value: 0  
Fairness Range: (-0.10, 0.10)



Ideal Value: 1.0  
Fairness Range: (0.80, 1.20)

A value  $< 1$  implies higher benefit for the privileged group and a value  $> 1$  implies a higher benefit for the unprivileged group.



Ideal Value: 0

Fairness is indicated by lower scores, higher scores represent inequality

# Age-based Bias is Made Evident in the German Credit Data Using the Metrics

---

- Comparing the metrics for bias based on Sex and Age
- Privileged Group: Male (Sex) and Old (Age > 25)
- Unprivileged Group: Female (Sex) and Young (Age < 25)

<i>Metric</i>	<i>Fairness Range</i>	<i>Sex</i>	<i>Age</i>
Statistical Parity Difference	(-0.10, 0.10)	0.01	-0.33
Equal Opportunity Difference	(-0.10, 0.10)	0.10	-0.42
Average Odds Difference	(-0.10, 0.10)	0.01	-0.31
Disparate Impact	(0.80, 1.20)	1.01	0.43
Theil Index	Lower the better	0.26	0.26

"original": *"He is a Nurse. She is a Optician. "* ("originalDistrib": [0.5, 0.5, 0.0])

Middle Language	Google	Yandex
tu *  Gender distinction lost or switched.	{..,"translated": "O hemşire. O bir Optisyendir.", "oto": "That nurse. It\u0026#39;s an Optic.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}	{..,"translated": "O bir Hemşire. Bir Gözlükçü.", "oto": "She\u0027s a nurse. An Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.5, 0.5]}
ru	{..,"translated": "Он медсестра. Она Оптик.", "oto": "He\u0026#39;s a nurse. She\u0026#39;s an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Он является медсестра. Она является Оптиком.", "oto": "He is a nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
it	{..,"translated": "Lui è un infermiere. Lei è un ottico.", "oto": "He is a nurse. She is an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Lui è un Infermiere. Lei è un Ottico.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
es	{..,"translated": "El es un enfermero. Ella es una Óptica.", "oto": "He is a nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Él es una Enfermera. Ella es un Oftalmólogo.", "oto": "He is a Nurse. She is an Ophthalmologist.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
hi *  Gender distinction replaced by both translators	{..,"translated": "वह नर्स है। वह एक ऑप्टिशियन है", "oto": "she\u0026#39;s a nurse. He is an optician", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "वह एक नर्स है. वह एक प्रकाशविज्ञानशास्त्री.", "oto": "She is a nurse. He is a optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
pt	{..,"translated": "Ele é um enfermeiro. Ela é uma óptica.", "oto": "He is a nurse. She\u0026#39;s an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Ele é uma Enfermeira. Ela é um Oculista.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
fr	{..,"translated": "Il est une infirmière. Elle est opticienne.", "oto": "He is a nurse. She is an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Il est une Infirmière. Elle est un Opticien.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
ar *  Gender distinction lost in Translation by both	{..,"translated": "هو نارس. وهي بصريات.", "oto": "It is Nars. They are optics.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}	{..,"translated": "هو ممرضة. هي العين.", "oto": "Is a nurse. Are the eyes.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}

# Mitigation Approaches

---

# Mitigation is a Socio-Technical Issue

---

- Removing problematic behavior – e.g., bias
  - **Concern**: do the developers understand the social implication of the original bias, and of any inserted by the remediation ?
  - **Concern**: what are the legal implications?
- Communicating behavior
  - Explaining decisions and characteristics via fairness metrics
    - **Concern**: which metric to use?
  - Third party evaluation and reproducible characterization of behavior on a scale
    - **Motivation**: nutrition labels in packaged food

# Bias Mitigation Algorithms Try to Improve the Fairness Metrics by Modifying Data, Model, or Predictions

---

The algorithms can be classified based on when a user can intervene in the machine learning pipeline:

## **Pre-processing (Data)**

- Reweighting
- Optimized Preprocessing

## **In-processing (Model)**

- Adversarial Debiasing

## **Post-processing (Predictions)**

- Reject Option Based Classification

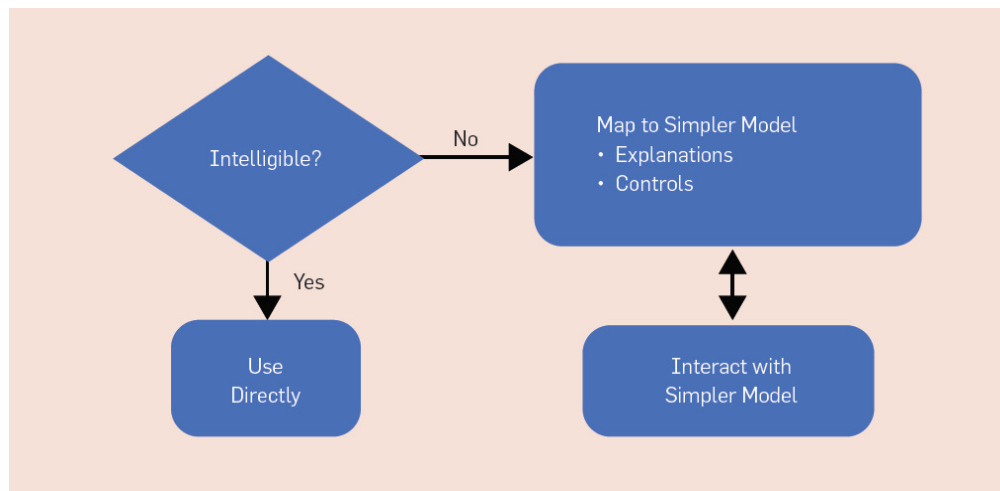
## **Guidance:**

*What type of mitigation to use depends on what stage the user can modify. Doing mitigation at the earliest is advisable.*

# Mitigation by Explanation

---

# Setting and Terminology: Intelligible Models and Explanations



- Transparency: providing stakeholders with relevant information about how a model works
- Explainability: Providing insights into model's behavior for specific datapoints

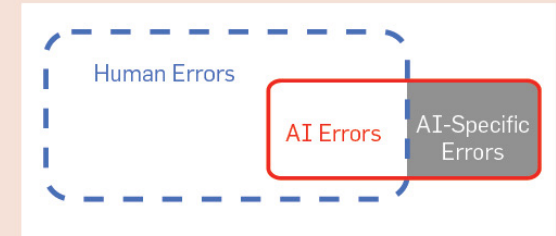
## Sources:

1. The Challenge of Crafting Intelligible Intelligence, Daniel S. Weld, Gagan Bansal, Communications of the ACM, June 2019, Vol. 62 No. 6, Pages 70-79, 10.1145/3282486
2. Explainable Machine Learning in Deployment, FAT\* 2020.



# Need for Intelligibility

The red shape denotes the AI's mistakes; its smaller size indicates a net reduction in the number of errors. The gray region denotes AI-specific mistakes a human would never make. Despite reducing the total number of errors, a deployed model may create new areas of liability (gray), necessitating explanations.



- **AI may have the wrong objective:** is AI right for the right reasons?
- **AI may be using inadequate features:** understand modeling issues
- **Distributional drift:** detect when and why models are failing to generalize
- **Facilitating user control:** guiding what preferences to learn
- **User acceptance:** especially for costly actions
- **Improving human insight:** improve algorithm design
- **Legal imperatives**

**Source:** The Challenge of Crafting Intelligible Intelligence, Daniel S. Weld, Gagan Bansal, Communications of the ACM, June 2019, Vol. 62 No. 6, Pages 70-79, 10.1145/3282486

# Types of Explanations

---

- **Feature-based:** from the features of the data, which feature(s) were most important for given decision output
  - Example: For a loan, is it income or the person's age ?
- **Sample-based:** from data in training, which data points were important for given test point; helps understand sampling and its representation in wider population
  - Example: For a loan, what instances similar to the loan application would have gotten the loan ?
- **Counter-factual:** what-ifs – what do you change about the input to change the decision output
  - Example: For a loan, does getting an additional borrower insurance increase chance of getting the loan?
- Natural language

**Source:** Explainable Machine Learning in Deployment, FAT\* 2020

# Stakeholders for Explanations

---

- Executives
  - Explainability as a market differentiator. Do we need explanations?
- ML engineers
  - How to improve model's performance?
- End-users
  - Understand business decisions emanating from usage of AI
    - Why was my load denied?
    - Why a particular treatment was recommended or de-prioritized ?
- Regulators
  - Prove that you did not discriminate based on existing laws

**Source:** Explainable Machine Learning in Deployment, FAT\* 2020

# References for AI Explainability

---

## Papers

- The Challenge of Crafting Intelligible Intelligence, Daniel S. Weld, Gagan Bansal, Communications of the ACM, June 2019, Vol. 62 No. 6, Pages 70-79, 10.1145/3282486
- “Why Should I Trust You?” Explaining the Predictions of Any Classifier, Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, in ACM’s Conference on Knowledge Discovery and Data Mining, KDD2016; <https://homes.cs.washington.edu/~marcotcr/blog/lime/>, <https://www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanations-lime/>
- Explainable Machine Learning in Deployment, FAT\* 2020, <https://arxiv.org/pdf/1909.06342.pdf>; Video: <https://www.youtube.com/watch?v=Hofl4uwxtPA>
- **Tutorial:** XAI tutorial at AAAI 2020, <https://xaitutorial2020.github.io/>
- **Tool:** AIX 360
  - Tool: <https://aix360.mybluemix.net/>
  - Video: <https://www.youtube.com/watch?v=Yn4yduyoQh4>
  - Paper: <https://arxiv.org/abs/1909.03012>

# LIME — Local Interpretable Model-Agnostic Explanations

---

**Paper:** “Why Should I Trust You?” Explaining the Predictions of Any Classifier, Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, ACM’s Conference on Knowledge Discovery and Data Mining, KDD2016

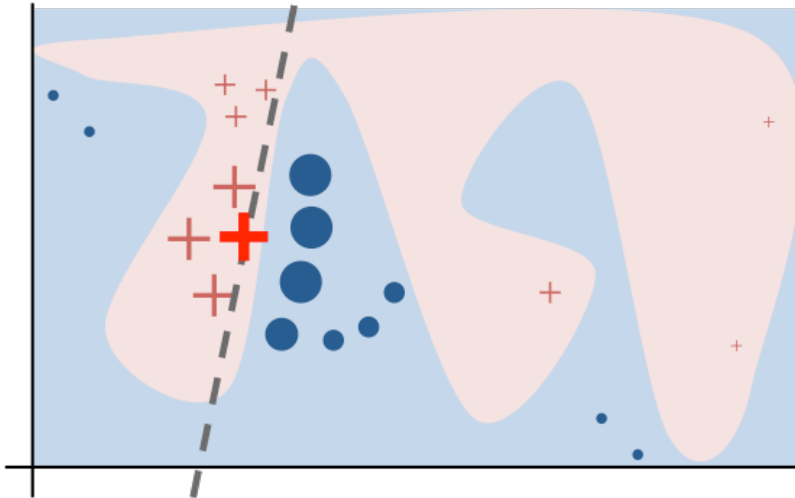
**Blogs:**

- <https://homes.cs.washington.edu/~marcotcr/blog/lime/>
- <https://www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanations-lime/>

**Code:** <https://github.com/marcotcr/lime>

# LIME Intuition

---



- Sample instances around X (point of interest)
- Weigh them according to their proximity to X
- Learn a linear model (dashed line) that approximates the model well in the vicinity of X
- Interpret the coefficients of the linear model based on data (features)

Source: <https://github.com/marcotcr/lime>

# LIME on Text

**Task:** predicting whether an is related to atheism (non-religions) or a particular religion (Christian)

**Question:** What is the classifier with >90% accuracy predicting based on ?

## Explanation:

“il we remove the words Host and NNTP from the document, we expect the classifier to predict atheism with probability  $0.58 - 0.14 - 0.11 = 0.31$ ”.

Prediction probabilities



atheism

christian



## Text with highlighted words

From: johnchad@triton.unm.edu (jchadwic)  
Subject: Another request for Darwin Fish  
Organization: University of New Mexico, Albuquerque  
Lines: 11  
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish.

This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.

Source: <https://github.com/marcotcr/lime>

# LIME on Image

**Question:** Why is this a frog?

Divide image into interpretable components - contiguous superpixels



Original Image



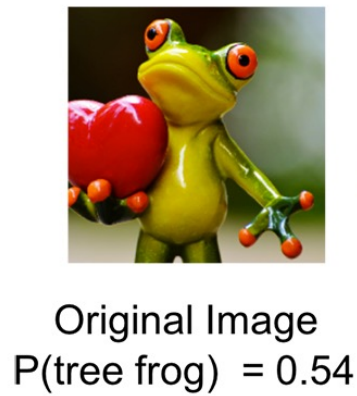
Interpretable  
Components





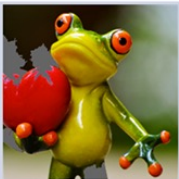

Source: <https://www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanations-lime/>

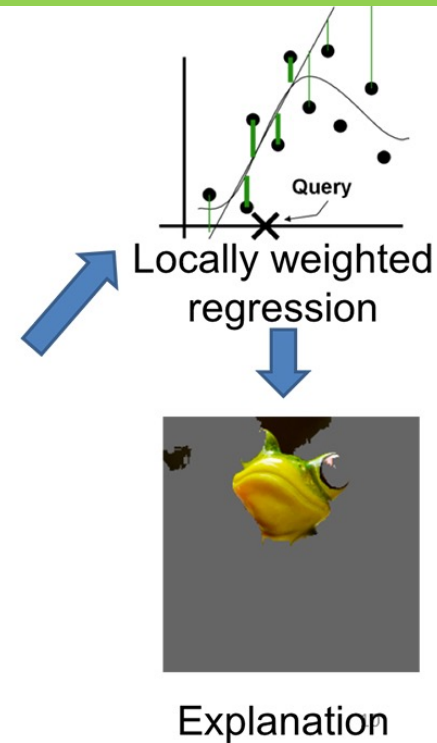


# LIME

1. Generate a data set of perturbed instances by turning some of the interpretable components “off” (gray).
2. For each perturbed instance, calculate probability that a tree frog is in the image according to the model.
3. Learn a simple (linear) model on this data set, which is locally weighted
4. Output regions with highest positive weights as an explanation, graying out everything else.



Perturbed Instances	$P(\text{tree frog})$
	 0.85
	 0.00001
	 0.52



# Another Method - ANCHOR

- Anchor – a rule that sufficiently describes the prediction locally such that changes to the rest of the feature values of the instance do not matter
- Example:  
<https://github.com/marcotcr/anchor/blob/master/notebooks/Anchor%20for%20text.ipynb>

**Paper:**  
<https://ojs.aaai.org/index.php/AAAI/article/view/1491>

**text = 'This is a good book .'**

**Anchor: good AND book AND is**  
**Precision: 0.97**

**Examples where anchor applies and model predicts positive:**

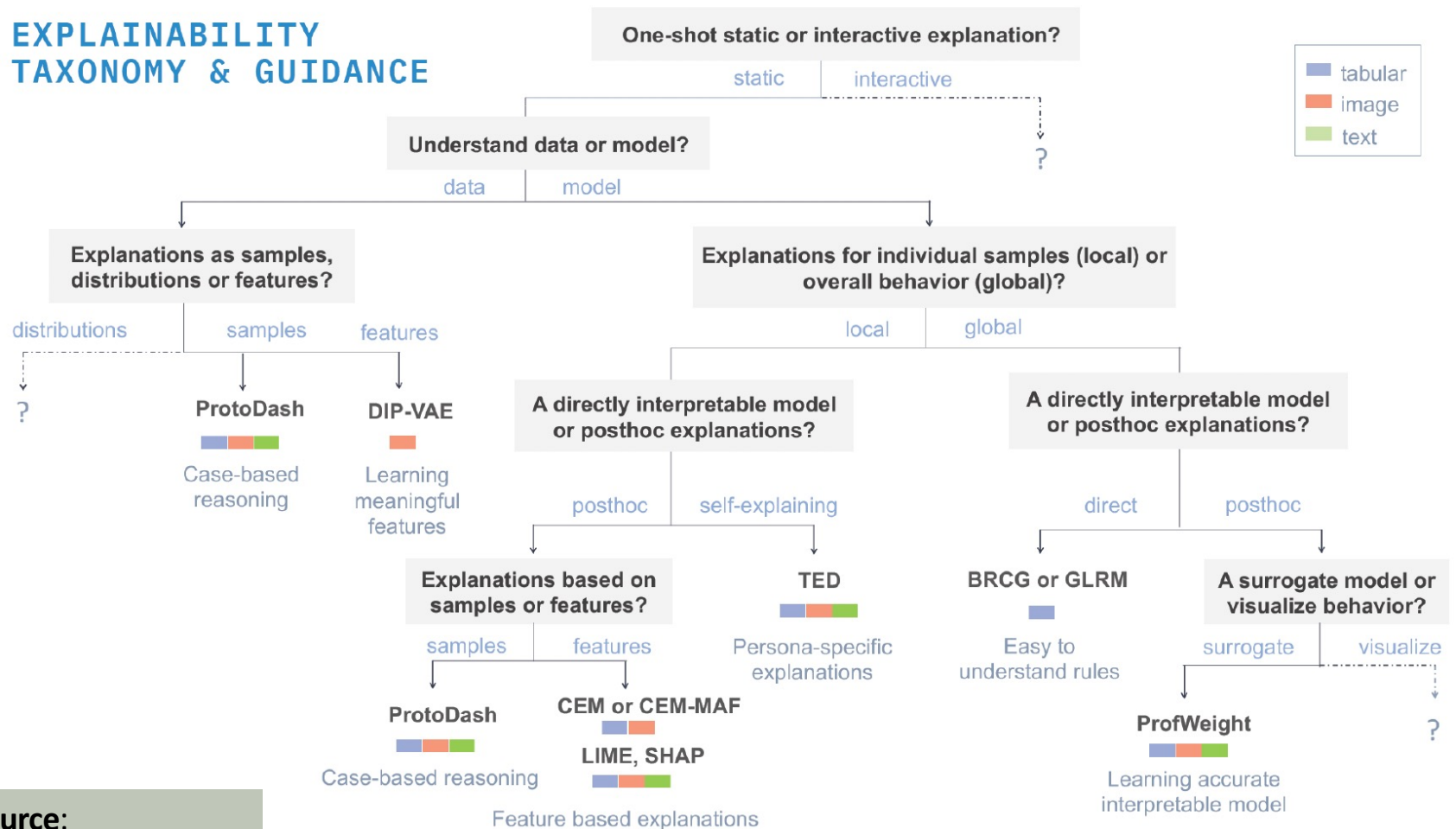
knowledge is a good book ; it is a good book . here is a good book :  
dawn is a good book ; another is a good book ; it is a good book .  
treasure is a good book . novels is a good book ; education is a good  
book . this is a good book .

**Examples where anchor applies and model predicts negative:**

everything is a good book . there is no good book . There is no good  
book here neither is a good book . nothing is another good book !

# A Spectrum of Explanations in AIX360

## EXPLAINABILITY TAXONOMY & GUIDANCE



**Slide Source:**  
Vera Liao's XAI talk 2020

# Mitigation by Rating for Trust

---

# Problem We Are Tackling for AI

## Insight

- Empower people to make informed decisions regarding which AI to choose
- Communicate trust information better!
  - Analogy: Food labels
- Facilitate users in understanding their choices

Calories 230		Calories from Fat 40
		% Daily Values*
Total Fat	8g	12%
Saturated Fat	1g	5%
Trans Fat	0g	
Cholesterol	0mg	0%
Sodium	160mg	7%
Total Carbohydrate	37g	12%
Dietary Fiber	4g	16%
Sugars	1g	
Protein	3g	
Vitamin A		10%
Vitamin C		8%
Calcium		20%
Iron		45%
Percent Daily Values are based on a diet of other people's secrets.		

Amount per 2/3 cup		Calories 230
		% DV*
12%	Total Fat 8g	
5%	Saturated Fat 1g	
	Trans Fat 0g	
0%	Cholesterol 0mg	
7%	Sodium 160mg	
12%	Total Carbs 37g	
14%	Dietary Fiber 4g	
	Sugars 1g	
	Added Sugars 0g	
	Protein 3g	
10%	Vitamin D 2mcg	
20%	Calcium 260mg	

In a series of previous work, we have developed ideas for rating bias of AI services

- For transactional services, method relies on a novel 2-stage testing method for bias. Papers in AIES 2018, IBM Sys Jour 2019, AAAI 2021 (Demo), IEEE Internet Computing 2021
- For conversation services (chatbot), method relies on testing properties (called issues) such as fairness, lack of information leakage, lack of abusive language, and adequate conversation complexity. Paper in IEEE Transactions on Technology and Society 2020.

But ideas are general and can apply to audio-, image- and multimodal AI services. Working on a generalized causal framework for rating

"original": *"He is a Nurse. She is a Optician."* ("originalDistrib": [0.5, 0.5, 0.0])

Middle Language	Google	Yandex
tu *  Gender distinction lost or switched.	{..,"translated": "O hemşire. O bir Optisyendir.", "oto": "That nurse. It\u0026#39;s an Optic.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}	{..,"translated": "O bir Hemşire. Bir Gözlükçü.", "oto": "She\u0027s a nurse. An Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.5, 0.5]}
ru	{..,"translated": "Он медсестра. Она Оптик.", "oto": "He\u0026#39;s a nurse. She\u0026#39;s an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Он является медсестра. Она является Оптиком.", "oto": "He is a nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
it	{..,"translated": "Lui è un infermiere. Lei è un ottico.", "oto": "He is a nurse. She is an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Lui è un Infermiere. Lei è un Ottico.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
es	{..,"translated": "El es un enfermero. Ella es una Óptica.", "oto": "He is a nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Él es una Enfermera. Ella es un Oftalmólogo.", "oto": "He is a Nurse. She is an Ophthalmologist.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
hi *  Gender distinction replaced by both translators	{..,"translated": "वह नर्स है। वह एक ऑप्टिशियन है", "oto": "she\u0026#39;s a nurse. He is an optician", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "वह एक नर्स है. वह एक प्रकाशविज्ञानशास्त्री.", "oto": "She is a nurse. He is a optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
pt	{..,"translated": "Ele é um enfermeiro. Ela é uma óptica.", "oto": "He is a nurse. She\u0026#39;s an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Ele é uma Enfermeira. Ela é um Oculista.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
fr	{..,"translated": "Il est une infirmière. Elle est opticienne.", "oto": "He is a nurse. She is an optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}	{..,"translated": "Il est une Infirmière. Elle est un Opticien.", "oto": "He is a Nurse. She is an Optician.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.5, 0.5, 0.0]}
ar *  Gender distinction lost in Translation by both	{..,"translated": "هو نارس. وهي بصريات.", "oto": "It is Nars. They are optics.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}	{..,"translated": "هو ممرضة. هي العين.", "oto": "Is a nurse. Are the eyes.", "values": ["He", "She", "OTHER"], "otoDistrib": [0.0, 0.0, 1.0]}

# Illustration of Translation Errors and Their Bias Perception About Gender

## Type: 1

**File:** Result - true - i - es.json

**Comment:** Gender changed to new value (It).

"**original**": "He is a Architect. She is a Astronomer.", "**translated**": "Es un arquitecto. Ella es una astrónoma.",

"**oto**": "It is an architect. She is an astronomer.",

## Type: 2

**File:** Result - false- g - tr.json;

**Comment:** Gender Flipped from She to He.

"**original**": "She is a Architect. She is a Astronomer.",

"**translated**": "O bir Mimar. O bir gökbilimcidir.",

"**oto**": "He is an architect. He is an astronomer.",

## Type: 5

**File:** \* - i - fr.json

**Comment:** No gender changed; Sentence modified.

"**original**": "He is a Accountant. He is a Actor /Actress.",

"**translated**": "C'est un comptable. C'est un acteur et un acteur.",

"**oto**": "He is an accountant. He is an actor and an actor.",

## Type: 3

**File:** Result - false- g - ar.json;

**Comment:** Grammatically wrong sentence; Subject missing.

"**original**": "She is a Factory worker. He is a Farmer.",

"**translated**": "هي عامل مصنع. هو مزارع.",

"**oto**": "~~Is~~ a factory worker. He is a farmer.",

## Type: 4

**File:** Result - false- g - tr.json;

**Comment:** "Multiple. Gender changed and flipped."

"**original**": "He is a Nurse. He is a Optician.",

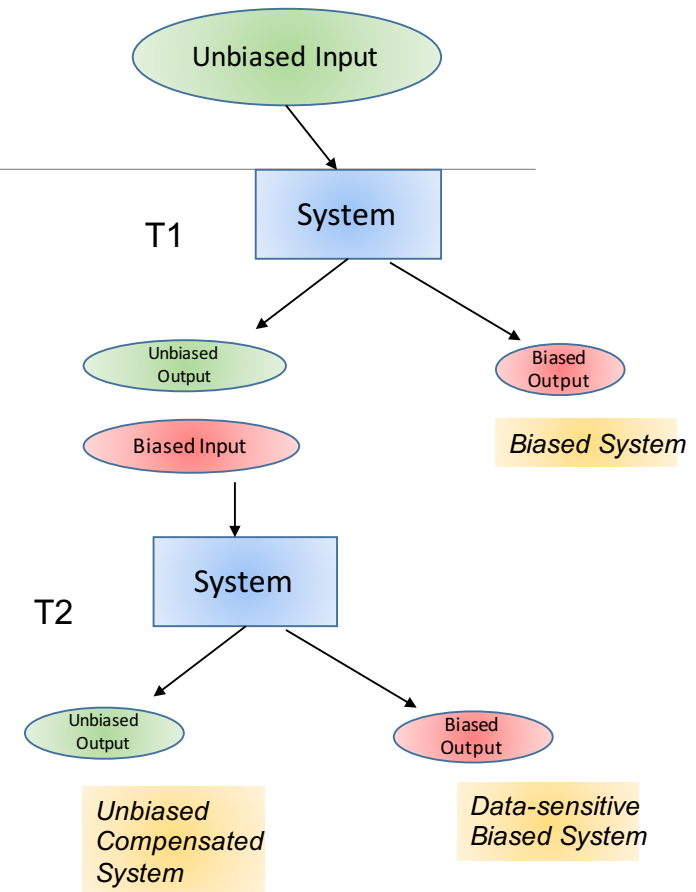
"**translated**": "O bir hemşire. O bir Optisyendir.",

"**oto**": "She is a nurse. It is an Optic.",

**1, 2, 3 and 4 have gender issues;  
3 and 5 have translation mistakes**

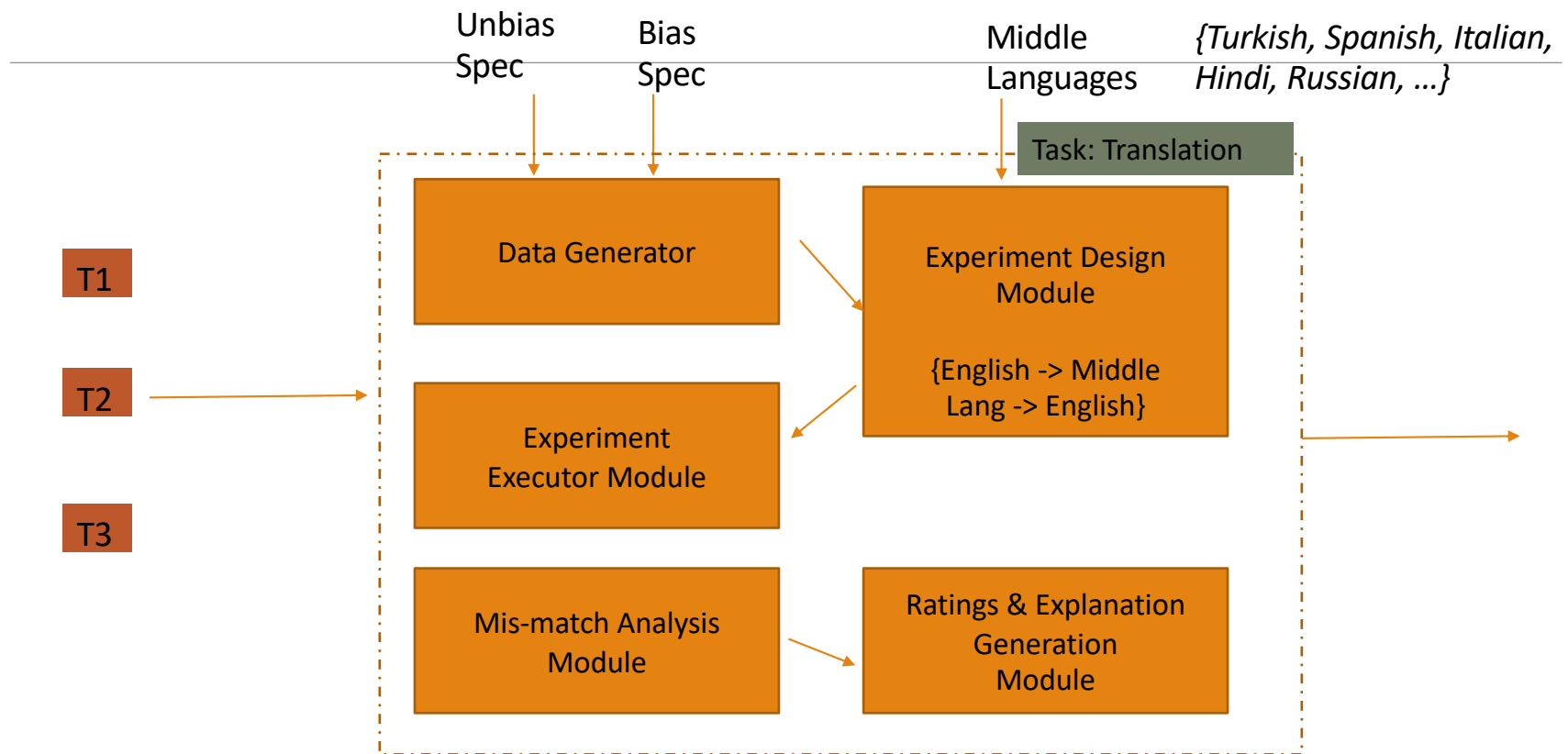
# Rating Translators

- We have an approach of 3rd party rating service: independent of API producer or consumer.
- Gives API producer distributions of biased and unbiased data.
- Does a new 2-step testing and produces ratings of 3 main levels: -
  - Unbiased Compensated System (UCS): Forces an assumed distribution among legal choices
  - Data-sensitive Biased System (DSBS): Its output follows a distribution similar to input
  - Biased System (BS): Follows a distribution statistically different from assumption
- Ratings supports multiple distribution definitions under unbiased and biased categories.
- Enhance scheme for compositions of APIs with their 3-level ratings
- Implementation and experiments on off-the-shelf translators and translation task with many middle languages.





# Illustrative Setup and Experiments



# But How Do People Perceive Ratings ? - VEGA Environment

---

Video: <https://www.youtube.com/watch?v=xZJklaRx4rQ>

Try the tool at: <http://vega-live.mybluemix.net/>

- Mariana Bernagozzi, Biplav Srivastava, Francesca Rossi and Sheema Usmani, VEGA: a Virtual Environment for Exploring Gender Bias vs. Accuracy Trade-offs in AI Translation Services, **AAAI 2021**. [Visualizing Ethics Rating, *Demonstration paper*]
- Mariana Bernagozzi, Biplav Srivastava, Francesca Rossi and Sheema Usmani, Gender Bias in Online Language Translators: Visualization, Human Perception, and Bias/Accuracy Trade-offs, **IEEE Internet Computing, Special Issue on Sociotechnical Perspectives**, Nov/Dec 2021 [Visualizing Ethics Rating, User Survey]

# Lecture 24: Concluding Comments

---

- The issue of trust
- Ethics and fairness issues
- Mitigation
  - Explanation and interpretation
  - Rating for trust
- For more details: Trusted AI course - <https://sites.google.com/site/biplavsrivastava/teaching/csce-590-trusted-ai>

# About Next Lecture – Lecture 25

---

# Lecture 25 Outline

---

- Use it to do Quiz 3
- Post questions on Piazza; discuss on Blackboard